

The Value of Information in Shortest Path Optimization

by

Michael David Rinehart

S.M., Massachusetts Institute of Technology (2005)

B.S., University of Maryland, College Park (2003)

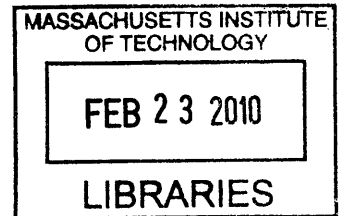
Submitted to the Department of Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

February 2010



© Massachusetts Institute of Technology 2010. All rights reserved.

ARCHIVES

Author
Department of Electrical Engineering and Computer Science
December 2, 2009

Certified by.....
Munther A. Dahleh
Professor
Thesis Supervisor

Accepted by.....
Terry P. Orlando
Chairman, Department Committee on Graduate Theses

The Value of Information in Shortest Path Optimization

by

Michael David Rinehart

Submitted to the Department of Electrical Engineering and Computer Science
on December 2, 2009, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

Abstract

Information about a random event (termed the source) is typically treated as a (possibly noisy) function of that event. Information has a destination, an agent, that uses the information to make a decision. In traditional communication systems design, it is usually assumed that the agent uses the information to produce an estimate of the source, and that estimate is in turn used to make the decision. Consequently, the typical objective of communication-systems design is to construct the communication system so that the joint distribution between the source and the information is “optimal” in the sense that it minimizes the average error of the estimate. Due to resource limitations such as cost, power, or time, estimation quality is constrained in the sense that the set of allowable joint distribution is bounded in mutual information.

In the context of an agent using information to make decisions, however, such metrics may not be appropriate. In particular, the true value of information is determined by how it impacts the average payoff of the agent’s decisions, not its estimation accuracy. To this end, mutual information may not be the most convenient measure of information quantity since its relationship to decision quality may be very complicated, making it difficult to develop algorithms for information optimization.

In this thesis, we study the value of information in an instance of an uncertain decision framework: shortest path optimization on a graph with random edge weights. Specifically, we consider an agent that seeks to traverse the shortest path of a graph subject to some side information it receives about the edge weights in advance of and during its travel. In this setting, decision quality is determined by the average length of the paths the agent chooses, not how often the agent decodes the optimal path.

For this application, we define and quantify a notion of information that is compatible with this problem, bound the performance of the agent subject to a bound on the amount of information available to it, study the impact of spreading information sequentially over partial decisions, and provide algorithms for information optimization. Meaningful, analytic performance bounds and practical algorithms for information optimization are obtained by leveraging a new type of geometric graph reduction for shortest path optimization as well as an abstraction of the geometry of sequential decision making.

Thesis Supervisor: Munther A. Dahleh

Title: Professor

Acknowledgments

First and foremost among those I owe my gratitude is my adviser Prof. Munther Dahleh. I want to first thank him for giving me the opportunity to join his research group and LIDS, the reasons for which of course go without saying. I also wish to thank him for the many lessons I learned from him while a student in his research group, and specifically for what I regard to be the most important of these lessons: how to seek out interesting questions and new landscapes for exploration. While I still have a few more years to go in refining this skill, I am fortunate to start from a such a strong base. Finally, I want to thank him for his excellent teaching and patience, which have contributed significantly to the very positive experience I had while a member of his group.

Within LIDS, I would like to thank some of the current and former students who have made my time here so enjoyable, specifically AmirAli Ahmadi, Aliaa Atwi, Ola Ayaso, Rose Faghih, Sleiman Itani, Sertac Karaman, Yola Katsargyri, Georgios Kotsalis, Patrick Kreidl, Nuno Martins, Sidhant Misra, Paul Njoroge, Mesrob Ohannessian, Mitra Osqui, Mardavij Roozbehani, Keith Santarelli, Sri Sarma, Pari Shah, Noah Stein, Danielle Tarraf, and Ermin Wei. Among these individuals, I want to specifically thank Yola for the many fun conversations, interesting stories, and good times; Paul for the many conversations and good jokes over free CSAIL tea and Muddy Charles beers; and my office mates (Yola, Mitra, and Rose) in part for their willingness to listen to my stories and far-out ideas, but mostly for just for making the office such a fun place to work.

Some friends outside of LIDS that I wish to thank are Al Avestruz, Jose Bohorquez, Wilfried Hofstetter, Deborah Howell, Behnam Jafarpour, Kelly Klima, Bill Nadir, and Crystal Ng. In particular, I would like to thank Al for countless enlightening conversations, good ideas, and just plain good times; Wilfried for many great discussions of geopolitics over many great beers; and Kelly for all the fun times while hanging out at the Thirsty. I also want to thank Lisa Vora and Jason Wampler for the good times and their hospitality whenever I visited Maryland, and I want to thank the many fun and interesting people I met while working with the Thirsty Ear Pub, including the staff (current and former), the executive committee, and the management. My time working with the Thirsty Ear holds a special place in my memories, and everyone involved with the pub is a part of those memories.

Finally, I wish to offer a special set of acknowledgments to my family and my girlfriend.

To my sister Kristin and my brothers Jason, Scott, and Tim, you make going back home more fun than any other destination that I can imagine, and you are collectively funnier than any other crowd I've met. You make traveling back to New Jersey a delight and traveling away a chore.

To my girlfriend Leila, thank you for your love, friendship, personality, sense of humor, support, unending patience, support, and unending patience. I have no doubt that you knew right away that wasn't a typo. There is also no doubt that I would not be here today if it was not for your support of my academic work and aspirations. At the same time, you keep me grounded and well balanced, and you keep me level when stress would otherwise have its own way. You are truly a rock for me as well as my best friend.

And to my parents, I want to thank you for your unending support of my goals at every step of the way. In particular, I want to thank my father for fostering and supporting the development of my interests, skills, and (by example) strong work ethic. I want to thank my mother for her constant belief in me (you were right!) and for her seeming unending wisdom that has had such a profound effect on me. For these reasons and more, I dedicate this dissertation to you both.

Michael Rinehart
November 24, 2009

This work was sponsored by the Army Research Office (W911NF-07-1-0568).

Contents

1	Introduction	10
2	Definitions, Notation, and Objectives	12
2.1	Basic Definitions and Notation	12
2.1.1	Random Variables and Sets	12
2.1.2	Graphs and Paths	13
2.2	Partial Information in Stochastic Optimization	14
2.2.1	Performance Formulation	14
2.2.2	Quantifying Information and Information Optimization	15
2.3	Partial Information in Shortest Path Optimization	16
2.3.1	Optional: Shortest Path Optimization under a Mutual Information Bound	17
2.3.2	Specializing to Gaussian Edge Weights	17
2.3.3	The Impacts of Graph Topology and Information Selection: Examples	18
2.4	Sequential Information in Shortest Path Optimization	20
2.4.1	The n -Stage Framework	21
2.4.2	Sequential Information Framework: Information Constraints	22
2.5	Objectives	22
2.6	Chapter Summary	22
3	Related Research	23
3.1	Order Statistics and Optimization	23
3.1.1	Order Statistics	23
3.1.2	Analytic Bounds based in Order Statistics	24
3.1.3	Computational Lower Bounds via the Generalized Chebyshev's In- equality	25
3.2	Stochastic Optimization	27
3.3	Shortest Paths on Random Graphs	27
3.4	Chapter Summary	28
4	The Value of Side Information in Shortest Path Optimization	29
4.1	Properties of the Path Polytope	29
4.1.1	Projection Matrix for the Path Polytope	29

4.1.2	Outer Spheric Approximation	30
4.1.3	Inner Spheric Approximation	31
4.2	Information Optimization using Graph Reductions	32
4.2.1	Information Optimization via Upper and Lower Bound Optimization	32
4.2.2	Upper Bound Optimization in the Gaussian Case	33
4.2.3	Lower Bound Optimization in the Gaussian Case	34
4.2.4	Special Case Analytic Solution for Information Optimization in the Gaussian Case	35
4.3	An Analytic Relationship between Capacity and Performance	37
4.3.1	Performance Lower Bounds	37
4.3.2	Performance Upper Bounds	38
4.4	Examples	40
4.4.1	Analytic Examples	40
4.4.2	Comparative Examples	42
4.5	Chapter Summary	45
5	The Value of Sequential Information in Shortest Path Optimization	46
5.1	Impact of Applying a Simple Information Constraint Set	46
5.2	The Uncontrolled Sequential Information Framework	48
5.2.1	Structuring Information	48
5.2.2	The Geometry of Partial Decisions in the $n = 2$ Case	49
5.2.3	Outer Approximating Partial Decisions in the $n = 2$ Case	50
5.2.4	Outer Approximating Partial Decisions in the General n Case	52
5.2.5	Analytic Performance Lower Bound under Uncontrolled Information	54
5.3	The Controlled Sequential Information Framework	58
5.3.1	Structuring Information	58
5.3.2	Performance Bounds	60
5.3.3	Examples	62
5.4	Chapter Summary	64
6	The Value of Information in Network Flow Optimization	65
6.1	Partial Information in Stochastic Social Welfare Optimization	65
6.1.1	Performance	65
6.1.2	Quantifying Information and Information Optimization	66
6.2	Network Flow Optimization under Limited Information	66
6.2.1	Specialization to Network Flow Optimization	66
6.2.2	Information Constraints	67
6.3	Outer Approximating Network Flow Optimization	67
6.4	Information Optimization via Lower Bound Optimization	69
6.4.1	The Impact of Information Selection: Examples	70
6.4.2	Computational Performance Bounds	72
6.4.3	Information Optimization	74
6.5	An Analytic Relationship between Capacity and Performance	75

6.6	Summary	76
7	Conclusions and Future Work	77
A	Additional Proofs and Results	79
B	The Value of Information in Energy Production	84
C	Upper Bounds via Splitting and Pruning	86
C.1	Computing an Upper Bound by Splitting the Graph	86
C.2	Upper Bound on Lost Performance from Pruning the Graph	87
C.2.1	Efficient Pruning in the Gaussian Case	88
C.2.2	Performance Loss from Path Pruning	90

List of Figures

2-1	Graph for Examples 1 and 2	18
2-2	Graph for Example 3.	19
4-1	$\Theta(x)$ (solid line) compared to $\min \{x, 0\}$ (dashed line).	38
4-2	The analytic bound of Theorem 7 compared to the optimization-based bound of Corollary 7 for a two-path graph. The solid lines are the analytic bound performances, and the asterisks are the optimization-based bound performances. Each line represents a graph with an increasing number of links per path.	44
4-3	The analytic bound of Theorem 7 compared to the optimization-based bound of Corollary 7 for random DAGs. The solid lines are the analytic bound performances, and the asterisks are the optimization-based bound performances. Each line represents a different random graph topology with the number of vertices fixed.	44
5-1	Example of a reduction in future decisions based on a past decision.	49
5-2	Illustration of a 3-dimensional ball with radius r intersecting an affine subspace S of dimension 2. The intersection is a 2-dimensional ball centered at c with radius r' (it appears as an ellipse due to distortion).	51
5-3	Illustration of a 3-dimensional ball intersecting an affine subspace S_1 of dimension 2 and then with another affine subspace $S_2 \subset S_1$ of dimension 1. With each intersection, we get a lower-dimensional ball (a line segment is a ball of dimension 1) with a new radius and center point.	52
5-4	Graph topology where sequential information is worse than non-sequential information.	62
5-5	Graph topology where (controlled) sequential information can achieve the same performance as the non-sequential case.	62
5-6	Graph topology where controlled sequential information can outperform non-sequential information.	63
6-1	Graph for Example 15	70
6-2	Graph for Example 16	71

6-3 Simulated performance of continuous flow with $R = 1$ on the graph in Figure 6-1 as capacity is varied from $0 \leq C \leq 20$. Each point represents a data point from the simulation. The solid line is a linear regression of those points. 76

Chapter 1

Introduction

Information about a random event (termed the *source*) is typically treated as a (possibly noisy) function of that event. Formally, information is as any other random event sharing a joint distribution with the source. Information is considered to have a destination, an *agent*, that uses the information to make a decision. In traditional communication systems design, it is usually assumed that the agent uses the information to produce an estimate of the source, and that estimate is in turn used to make the decision. Consequently, the typical objective of communication systems design is to construct the communication system so that the joint distribution between the source and the information is “optimal” in the sense that it minimizes the average error of the estimate. Due to resource limitations such as cost, power, or time, information optimization is constrained in the sense that the joint distribution is limited to a set of distributions.

A fundamental framework for the study of information communication and estimation error is information theory. In information theory, the total “amount” of information contained in the source is defined as its *entropy*, and the *mutual information* between the source and the information is defined as the amount by which the entropy of the source is reduced on average once we know the information. Bounds on the agent’s estimation error can be derived or computed directly from the mutual information. The limitations imposed on optimizing the joint distribution between the source and the information is captured in the abstraction of a *communication channel* over which information is communicated. The communication channel limits the maximum possible mutual information between the source and the agent. This bound is termed the *channel’s capacity*.

In the context of an agent using information to make decisions, such traditional metrics may not be appropriate. In particular, the true value of information is determined by how it impacts the average payoff of the agent’s decisions, not its estimation quality. To this end, mutual information may not be the most convenient measure of information quantity since its relationship to decision quality may be very complicated.

In this thesis, we study the value of information in an uncertain decision framework. We consider an agent who (a) uses information about the source to estimate the quality of its possible decisions (not necessarily the source), (b) can optimize the information it receives

subject to a bound on some useful measure for information quantity, and (c) can consider making partial decisions as information arrives or wait until all of the information arrives before making a complete decision. We specifically explore this framework in the context of shortest path optimization.

In shortest path optimization, an agent uses information about the random edge weights of a graph to determine the shortest-average path in the graph. In this setting, information quality is determined by the average length of the paths the agent chooses, not how often the agent decodes the optimal path. We define and quantify a notion of information that is compatible with this problem, bound the performance of the agent subject to a bound on the amount of information available to it, study the impact of spreading information sequentially over partial decisions, and provide algorithms for optimizing information. Meaningful, analytic performance bounds and practical algorithms for information optimization are obtained by leveraging a new type of geometric graph reduction for shortest path optimization as well as an abstraction partial decision making.

The outline of this thesis is as follows. In Chapter 2, we define the concepts and notations used throughout this thesis, formally describe the framework for shortest path optimization under limited information, and present our objectives. In Chapter 3, we review some research that is, in some cases, qualitatively related to our work and, in other cases, technically related. In Chapter 4, we provide computational and analytic bounds for performance and algorithms for information optimization. In Chapter 5, we study how information spread over partial decisions impacts the agent’s performance. In Chapter 6, we study of an alternative but related problem: network flow optimization. And, finally, in Chapter 7, we provide some final conclusions.

Chapter 2

Definitions, Notation, and Objectives

The concepts and tools used in this thesis are based in probability theory, optimization theory, and graph theory, and a working knowledge of these fields is assumed of the reader. In this chapter, we present a set of notation, assumptions, and basic definitions that will be used throughout the remainder of the thesis. A limited set of foundational results for the thesis are also provided.

2.1 Basic Definitions and Notation

2.1.1 Random Variables and Sets

We write random variables (RVs) in capital letters (e.g., X), and denote the event $X \in S$ as $P(X \in S)$ or $P_X(S)$. We write $X \sim p$ if X has p as its probability density function (PDF), and let $N(\mu, \sigma^2)$ be the normal distribution with mean μ and variance σ^2 . $E[X]$ and $\text{VAR}[X]$ are, respectively, the expected value and variance of X , and for a random vector $X = (X_1, \dots, X_n)$, $\text{VAR}[X] = \sum_{i=1}^n \text{VAR}[X_i]$ whereas $\text{COV}[X] = E[XX^T] - E[X]E[X]^T$. For two RVs X, Y , we define $\hat{X}(Y) = E[X|Y]$ as the estimate of X given Y (which we simplify to \hat{X} if the argument is understood), and we say $X \stackrel{d}{=} Y$ if both RVs are drawn from the same distribution.

If A is a set, $|A|$ is the number of elements in A . For another set B , $A \setminus B$ is the set of elements in A but not in B . If $A \subset \mathbb{R}^n$ and $x \in \mathbb{R}^n$, $A - x = \{a - x \mid a \in A\}$. If A is a subspace, then A^\perp is the orthogonal subspace to A , and if A is an affine subspace, then define the subspace $\dot{A} = A - \{a\}$ for any $a \in A$ ¹. We denote the sphere of radius r and center c as $B(r, c) = \{x \mid \|x - c\|^2 \leq r^2\}$.

For a positive (semi)definite matrix M , $N = \sqrt{M}$ is the unique positive (semi)definite matrix satisfying $M = N^2$.

¹It does not matter which $a \in A$ is selected, \dot{A} will be the same subspace

Finally, for $a \leq b$, let

$$[x]_a^b = \min \left\{ \max \{x, a\}, b \right\},$$

with $[x]^+ = \max \{x, 0\}$, and $[x]^- = \min \{x, 0\}$. For ease, we simply denote $E \left[[x]_a^b \right]$ as $E [x]_a^b$.

2.1.2 Graphs and Paths

We define a graph G by a pair (V, E) of vertices V and edges E . Because we allow any two vertices to have multiple edges connect them, we forgo the usual definition $E \subset V \times V$ and instead define a head and tail for each edge $e \in E$ by $\text{hd}(e) \in V$ and $\text{tl}(e) \in V$ respectively.

Each edge e in the graph is associated with an *edge weight* w_e . The vector of all weights is $w = [w_1 \dots w_{|E|}]^T$. Because we consider edge weights to be random, we write the vector as W , and we assume that the probability distribution is known. Finally, we denote the first and second moments of W by

$$\begin{aligned} \mu &= E[W], \mu_e = E[W_e], \\ \Lambda_W &= \text{COV}[W], \sigma_e^2 = \text{VAR}[W_e]. \end{aligned}$$

We now define the notion of a *path* in the graph.

Definition 1 (Path). *A sequence $p = (e_1, e_2, \dots, e_n)$ of edges is a path if $\text{tl}(e_i) = \text{hd}(e_{i+1})$, and we say p goes from $v_1 = \text{hd}(e_1)$ to $v_{n+1} = \text{tl}(e_n)$.*

Definition 2 (Acyclic Path). *A path $p = (e_i)$ is acyclic if there are no two indices $i < j$ such that $\text{hd}(e_i) = \text{tl}(e_j)$.*

Assumption 1 (DAG). *G is a directed acyclic graph (DAG) (i.e., all paths the p of G are acyclic).*

We also assume the existence of two vertices $s, t \in V$, respectively termed the *start* and *termination* vertices, that (uniquely) satisfy the following assumption.

Assumption 2. *There is a path from s to each vertex $v \in V \setminus \{s\}$ as well as a path from each vertex $v \in V \setminus \{t\}$ to t .*

Let $P = P(G)$ be the set of all paths from s to t in G . With some abuse of notation, we can write each $p \in P$ as a 0-1 vector in $\mathbb{R}^{|E|}$ where $p_e = 1$ if $e \in p$ and $p_e = 0$ otherwise. In this case, P is a set of all such vectors in $\mathbb{R}^{|E|}$. Using our vector notation, the length of a path $p \in P$ is simply $p^T W$.

Let $\mathcal{P} = \text{convex hull}\{P\}$. For an edge weight vector w , the length of the shortest path in the graph can be written as

$$\min_{p \in P} \{p^T w\} = \min_{p \in \mathcal{P}} \{p^T w\}.$$

Proposition 1.

$$\mathcal{P} = \{x \in [0, 1]^{|E|} \text{ such that } \sum_{e \in \text{tl}(v)} x_e - \sum_{e \in \text{hd}(v)} x_e = \begin{cases} 1 & v = s \\ -1 & v = t \\ 0 & \text{otherwise} \end{cases}\}$$

Proof. Let X be the set described in the theorem. It is well known that the shortest path under any set of edge weights w can be expressed as $\min_{p \in X} \{p^T w\}$. Therefore $X = \mathcal{P}$. \square

2.2 Partial Information in Stochastic Optimization

To motivate our framework for studying the value of side information in shortest path optimization, we present the basic ideas in a general stochastic optimization context. For the purposes of this section, let W be some RV with a known distribution.

2.2.1 Performance Formulation

Consider the following stochastic optimization:

$$J(W) = \min_{x \in \mathcal{X}} h(x, W).$$

Clearly, since W is a RV, $J(W)$ is also a RV, and so the average performance of the optimization is $\mathbb{E}[J(W)]$.

Consider now the task of finding an “optimal” decision x without having the realization of W . A reasonable objective is to select the x that minimizes the average of the objective:

$$\bar{J} = \min_{x \in \mathcal{X}} \mathbb{E}[h(x, W)].$$

Since \bar{J} is a constant, $\mathbb{E}[\bar{J}] = \bar{J}$. By Jensen’s Inequality, $\mathbb{E}[J(W)] \leq \bar{J}$.

We call the first case (where the realization of W was known) the *full-information* case. We call the latter case the *zero-information* case.

We are interested in formulating an in-between *partial-information* case. To this end, we introduce another RV Y that represents the agent’s *side information* about W and write the optimization as a function of our side information:

$$J(Y) = \min_{x \in \mathcal{X}} \mathbb{E}[h(x, W)|Y].$$

Once again, because Y is a RV, $J(Y)$ is also a RV, and so the average performance under Y is simply $\mathbb{E}[J(Y)]$. Clearly, the information Y contains about W is completely determined by their joint-distribution p_{WY} , so we define $J(p_{WY}) = \mathbb{E}[J(Y)]$:

$$J(p_{WY}) = \mathbb{E} \left[\min_{x \in \mathcal{X}} \mathbb{E}[h(x, W)|Y] \right]. \quad (2.1)$$

Remark 1. *Intuitively, we are "averaging-out" the information about $h(x, W)$ that we do not have from Y , much like in the zero-information case. The full- and zero-information cases are easily obtained by substituting $Y = 0$ (a constant) and $Y = W$ to respectively yield $J(Y) = \bar{J}$ and $J(Y) = J(W)$.*

Proposition 2. $E[J(W)] \leq J(p_{WY}) \leq \bar{J}$

Proof. The proof is a simple application of Jensen's Inequality:

$$\begin{aligned} E[J(W)] &= E \left[\min_{x \in X} h(x, W) \right] = E \left[E \left[\min_{x \in X} h(x, W) \mid Y \right] \right] \\ &\leq E \left[\min_{x \in X} E[h(x, W) \mid Y] \right] = J(p_{WY}) \\ &\leq \min_{x \in X} E[E[h(x, W) \mid Y]] = \min_{x \in X} E[h(x, W)] = \bar{J}. \end{aligned}$$

□

2.2.2 Quantifying Information and Information Optimization

The agent is also given some flexibility in determining the side information it receives in the form of being able to choose the joint distribution p_{WY} . Without added constraints, though, the agent will choose a distribution that yields $\hat{W}(Y) = W$. Therefore, we define a "bound" Γ to limit the set of allowable distributions. In terms of performance, we seek the solution to

$$J(\Gamma) = \min_{p_{WY} \in \Gamma} J(p_{WY}). \quad (2.2)$$

We call (2.2) the *information optimization*.

To quantify information, we further generalize the concept of an information bound to a family of constraint sets $\{\Gamma(C)\}$ parameterized by a non-negative scalar C called the *capacity*. For ease, we simplify our notation by writing

$$J(C) = J(\Gamma(C))$$

and call $J(C)$ the *optimal performance under capacity C* . Although not critical to the analysis of this thesis, desirable properties of $\Gamma(C)$ include:

- $\Gamma(C_1) \subset \Gamma(C_2)$ if $C_1 \leq C_2$,
- $p_{WY} \in \Gamma(0)$ implies $E[W \mid Y] = E[W]$,
- if $p_{W(Y_1 Y_2)} \in \Gamma(C)$, then p_{WY_1} and p_{WY_2} are in $\Gamma(C)$ ²,
- and there exists a C_{max} such that there exist a $p_{WY} \in \Gamma(C_{max})$ satisfying $\hat{W}(Y) = W$.

²Note that Y is any side information, and so it can be taken as a tuple of side information as well

2.3 Partial Information in Shortest Path Optimization

We now specialize our framework to shortest path optimization. We begin by defining the information constraint sets $\{\Gamma(C)\}$:

$$\Gamma(C) = \{p_{WY} \mid \text{VAR}[E[W|Y]] \leq C\}.$$

Our choice of $\Gamma(C)$ is a practical selection motivated by the analysis that is to follow in this paper. It furthermore obeys our desired properties for $\{\Gamma(C)\}$.

Proposition 3. $0 = \text{VAR}[E[W]] \leq \text{VAR}[E[W|Y_1]] \leq \text{VAR}[E[W|Y_1Y_2]] \leq \text{VAR}[W]$

Proof. The first equality is obvious. The remainder of the inequalities are based on the identity $E[f(Y)E[W|Y]] = E[f(Y)W]$. Without loss of generality, assume $W \in \Re$ and $E[W] = 0$. For any Y ,

$$\begin{aligned} 0 &\leq \text{VAR}[E[W|Y] - W] \\ &= \text{VAR}[E[W|Y]] + \text{VAR}[W] - 2E[E[W|Y]W] \\ &= \text{VAR}[E[W|Y]] + \text{VAR}[W] - 2E[\hat{W}(Y)W] \\ &= \text{VAR}[E[W|Y]] + \text{VAR}[W] - 2E[\hat{W}(Y)E[W|Y]] \\ &= \text{VAR}[E[W|Y]] + \text{VAR}[W] - 2\text{VAR}[E[W|Y]] \\ &= \text{VAR}[W] - \text{VAR}[E[W|Y]], \end{aligned}$$

The remaining inequalities are similarly proved. \square

The interpretation of Proposition 3 is that as we add information to our estimate in the form of $Y = (Y_1)$ to $Y = (Y_1, Y_2)$, our measure for information increases. The lower bound represents the case of having zero information, and the upper bound represents the case of having full information.

Now, given a joint distribution p_{WY} between the edge weights W of the graph and the information Y that the agent receives, we can write the agent's average performance as

$$J(p_{WY}) = E \left[\min_{p \in P} \{p^T E[W|Y]\} \right], \quad (2.3)$$

which simplifies to

$$J(p_{\hat{W}}) = E \left[\min_{p \in P} \{p^T \hat{W}\} \right]. \quad (2.4)$$

Notice that (2.4) only depends on $p_{\hat{W}}$ ³. We can also equivalently parameterize $\Gamma(C)$ by

$$\Gamma(C) = \left\{ p_{\hat{W}} \mid \text{VAR}[\hat{W}] \leq C \text{ and there is a } Y \text{ so that } \hat{W} \stackrel{d}{=} E[W|Y] \right\} \quad (2.5)$$

³This is true for any linear objective, not just shortest path optimization.

and ignore the joint distribution p_{WY} altogether.

2.3.1 Optional: Shortest Path Optimization under a Mutual Information Bound

Why do we define a family of abstract sets $\{\Gamma(C)\}$ as our information bound and not simply use mutual information? In general, we want to apply bounds $\Gamma(C)$ that yield a nice relationship between C and $J(C)$. For instance, we will see that our variance bounds relate for information relates nicely the shortest path optimization.

We can relate our information bound to mutual information, however. Let

$$\Gamma'(C) = \{p_{WY} | I(W; Y) \leq I_{\max}(C)\}$$

where $I_{\max}(C) = \max_{p_{WY} \in \Gamma(C)} \{I(W; Y)\}$. Then $\Gamma(C) \subset \Gamma'(C)$ so that $J(\Gamma'(C)) \leq J(C)$.

In general, computing $I_{\max}(C)$ is difficult, but so may be computing the performance under mutual information bounds directly. For instance, it is straightforward to see that the sets $\{\Gamma(C)\}$ that we are using for shortest path optimization relate to mutual information via a rate-distortion problem:

$$I_{\max}(C) = \min_{\{p_{WY} | \text{VAR}[W - \hat{W}] \leq \text{VAR}[W] - C\}} \{I(W; Y)\}.$$

In general, it is not trivial to solve this optimization, especially in the multivariable case.

2.3.2 Specializing to Gaussian Edge Weights

A particular subcase of interest to us is that of Gaussian edge weights. If $Y = W + N$ where W and N are independent, $W \sim N(\mu, \Lambda_W)$ with $\Lambda_W > 0$, and $N \sim N(0, \Lambda_N)$, then

$$\hat{W}(Y) = \Lambda_W (\Lambda_W + \Lambda_N)^{-1} (Y - \mu) + \mu.$$

Information optimization in this special case is equivalent to designing the distribution of the noise N , or, equivalently, its covariance matrix Λ_N . It is straightforward to show that designing a positive semidefinite Λ_N is equivalent to designing $\Lambda_{\hat{W}}$ (which is equal to $\Lambda_W (\Lambda_W + \Lambda_N)^{-1} \Lambda_W$) directly subject to several constraints. Denote this new constraint set as $\Gamma_G(C)$:

$$\Gamma_G(C) = \{\Lambda_{\hat{W}} \mid 0 \leq \Lambda_{\hat{W}} \leq \Lambda_W \text{ and } \text{Tr}(\Lambda_{\hat{W}}) \leq C\}. \quad (2.6)$$

Note that $\Gamma_G(C)$ is a convex set. It remains convex if we add additional convex constraints such as $\Lambda_{\hat{W}} \sim \text{diag}$.

Finally, for ease, denote $J_G(C) = J(\Gamma_G(C))$.

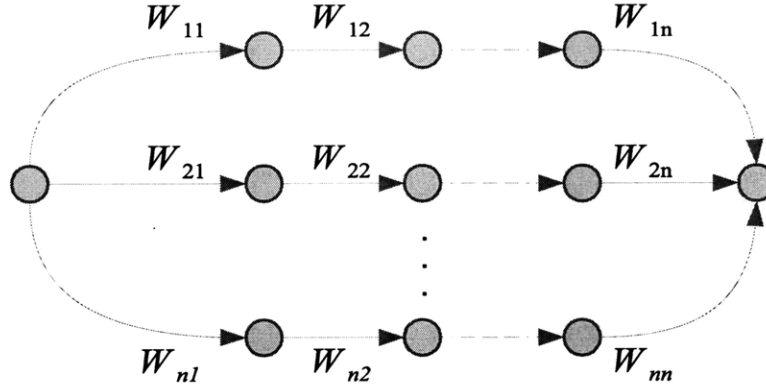


Figure 2-1: Graph for Examples 1 and 2

2.3.3 The Impacts of Graph Topology and Information Selection: Examples

Using our notation from the Gaussian case, we can highlight the impact that even a modest information optimization can have on performance.

Example 1. Let G be the graph in Figure 2-1 having n disjoint paths $\{p_1, \dots, p_n\}$ from s to t with each path having n edges each. Let $W_{ij} \sim N(0, 1)$ be the (random) weight of edge j on path i and assume the edge weights are independent. Let the “distribution” $\Lambda^s \in \Gamma_G(n)$ satisfy $\Lambda^s \sim \text{diagonal}$, $\Lambda_{ee}^s = 1$ if $e \in p_1$, and $\Lambda_{ee}^s = 0$ otherwise. Essentially, the estimates \hat{W} only contain information about the edges in p_1 , meaning that $\hat{W}_e = W_e$ for $e \in p_1$ and $\hat{W}_e = 0$ for $e \notin p_1$.

For this distribution, the average performance is

$$\begin{aligned}
 J(\Lambda^s) &= \mathbb{E} \left[\min_i \left\{ \sum_j \hat{W}_{ij} \right\} \right] \\
 &= \mathbb{E} \left[\min \left\{ \sum_j W_{1j}, 0 \right\} \right] \\
 &= \mathbb{E} \left[\min \{ \sqrt{n}Z, 0 \} \right] \\
 &= \sqrt{n} \mathbb{E} \left[\min \{ Z, 0 \} \right] \\
 &= -\frac{1}{\sqrt{2\pi}} \sqrt{n}
 \end{aligned}$$

where $\sum_j W_{1j} \stackrel{d}{=} \sqrt{n}Z$ with $Z \sim N(0, 1)$.

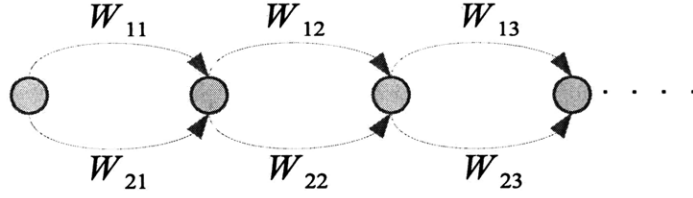


Figure 2-2: Graph for Example 3.

Example 2. Let G be the graph in Figure 2-1 and take the distribution $\Lambda^p \in \Gamma_G(n)$ where $\Lambda^p \sim \text{diagonal}$, $\Lambda_{ee}^p = 1$ if e is the first edge of any path p_i , and $\Lambda_{ee}^p = 0$ otherwise. Essentially, the estimates \hat{W} only contain information about the first edge in each path, meaning that $\hat{W}_e = W_e$ if e is one of these links and $\hat{W}_e = 0$ otherwise.

For this distribution, the average performance is

$$\begin{aligned}
 J(\Lambda^p) &= \mathbb{E} \left[\min_i \left\{ \sum_j \hat{W}_{ij} \right\} \right] \\
 &= \mathbb{E} \left[\min \{W_{11} + 0, W_{21} + 0, \dots, W_{n1} + 0\} \right] \\
 &= \mathbb{E} \left[\min \{Z_{11}, Z_{21}, \dots, Z_{n1}\} \right] \\
 &\geq -\sqrt{2 \ln n}
 \end{aligned}$$

where the last inequality is obtained by using Lemma 3 in [8].

There is a significant difference between the performances of the two examples. If we increase n , the average performance yielded from applying Λ^s outstrips that obtained using Λ^p quite substantially. This motivates our desire to optimize the information received by the agent.

The next example presents a topology for which the agent's performance significantly improves with capacity.

Example 3. Consider the graph in Figure 5-5 with $W_{ij} \sim N(0, 1)$ and independent and assume $C < \frac{|E|}{2}$. Choose any $\Lambda \in \Gamma_G(C) \cap \{\Lambda \text{ is diagonal}\}$, and denote $\lambda_e = \Lambda_{ee}$. We have

$$\begin{aligned}
 J(p_{\hat{W}}) &= \sum_j \mathbb{E} \left[\min \{ \hat{W}_{1j}, \hat{W}_{2j} \} \right] = \sum_j \mathbb{E} \left[\hat{W}_{2j} \right] + \mathbb{E} \left[\min \{ \hat{W}_{1j} - \hat{W}_{2j}, 0 \} \right] \\
 &= J(0) + \sum_j \min \left\{ \sqrt{\lambda_{1j}^2 + \lambda_{2j}^2} Z, 0 \right\} = J(0) - \frac{1}{\sqrt{2\pi}} \sum_j \lambda_j,
 \end{aligned}$$

where $\lambda_j = \sqrt{\lambda_{1j}^2 + \lambda_{2j}^2}$ and $Z \sim N(0, 1)$.

The optimal $\{\lambda_j^*\}$ (there are $|E|/2$ of them) are given by the solution to

$$\max_{\lambda} \left\{ \sum_j \lambda_j \right\} \text{ subject to } \sum_j \lambda_j^2 = C$$

yielding $\lambda_j^* = \sqrt{\frac{C}{|E|/2}}$. Thus,

$$J_G(C) \leq J(0) - \frac{1}{\sqrt{2\pi}} \frac{|E|}{2} \sqrt{\frac{C}{|E|/2}} = J(0) - \frac{1}{2\sqrt{\pi}} \sqrt{|E|} \sqrt{C}.$$

Example 4. Finally, we consider a path having a single path from s to t . In this case, we get

$$J(C) = \mathbb{E} \left[\sum_e \hat{W}_e \right] = \sum_e \mathbb{E} [\hat{W}_e] = J(0),$$

where $J(0)$ is the performance under no information.

2.4 Sequential Information in Shortest Path Optimization

We now consider the impact of spreading information sequentially over partial decisions, called the *sequential information* framework. This case is quite similar to the previous non-sequential case. We begin by developing the two-stage sequential information framework. The motivation for the formal framework is the following sequence of events:

1. The agent receives some side information Y_1 while it is at vertex s .
2. It estimates the edge weights of the graph using Y_1 .
3. It chooses a successor vertex v of s to which to travel.
4. After arriving at v , it receives some side information Y_2 .
5. It reestimates the edge weights using Y_1 and Y_2 .
6. Finally, it chooses the path starting from v with the smallest expected length.

Remark 2. A caveat in step 3 is that the agent does not necessarily choose the vertex v that falls along the shortest-path estimated using Y_1 . The agent may consider using knowledge of the fact that it will receive future information to influence its choice.

The side information in the two-stage framework is given by a tuple (Y_1, Y_2) that shares a joint distribution $p_{WY_1Y_2}$ with W . The agent's average performance in this two-stage

setup is

$$J(p_{WY_1Y_2}) = \mathbb{E} \left[\min_{e_1 | \text{hd}(e_1)=s} \left\{ \mathbb{E}[W_{e_1}|Y_1] + \mathbb{E} \left[\min_{\text{path } p=(e_1, \dots)} \left\{ \sum_{e_1 \neq e \in p} \mathbb{E}[W_e|Y_1Y_2] \right\} | Y_1 \right] \right\} \right].$$

When making the decision about the first edge e_1 to take, the agent only knows Y_1 as well as its average performance after e_1 subject to receiving more side information. For a given choice e_1 , the agent can only take paths p after e_1 (i.e., paths of the form $p = (e_1, \dots)$). It chooses the remainder of p by using its knowledge of both Y_1 and Y_2 .

To simplify notation, let $\hat{W}_i = \mathbb{E}[W|Y_i]$ and $\hat{W}_{12} = \mathbb{E}[W|Y_1Y_2]$. We can write the two-stage performance more compactly.

Proposition 4.

$$J(p_{WY_1Y_2}) = \mathbb{E} \left[\min_{e | \text{hd}(e)=s} \left\{ \mathbb{E} \left[\min_{\text{path } p \in P | p_e=1} \left\{ p^T \hat{W}_{12} \right\} | Y_1 \right] \right\} \right] \quad (2.7)$$

Proof. First, substitute our \hat{W}_i notation to get

$$J(p_{WY_1Y_2}) = \mathbb{E} \left[\min_{e_1 | \text{hd}(e_1)=s} \left\{ (\hat{W}_1)_{e_1} + \mathbb{E} \left[\min_{\text{path } p=(e_1, \dots)} \left\{ \sum_{e_1 \neq e \in p} (\hat{W}_{12})_e \right\} | Y_1 \right] \right\} \right]$$

Using (a) the relationship $\hat{W}_1 = \mathbb{E}[\hat{W}_1|Y_1] = \mathbb{E}[\mathbb{E}[\hat{W}_1|Y_1Y_2] | Y_1] = \mathbb{E}[\hat{W}_{12}|Y_1]$, and (b) that $(\hat{W}_1)_{e_1}$ is a constant with respect to the inner minimization, we can pull $(\hat{W}_1)_{e_1} = \mathbb{E}[(\hat{W}_{12})_{e_1}|Y_1]$ into the inner expectation and minimization to get

$$= \mathbb{E} \left[\min_{e_1 | \text{hd}(e_1)=s} \left\{ \mathbb{E} \left[\min_{\text{path } p=(e_1, \dots)} \left\{ \sum_{e \in p} (\hat{W}_{12})_e \right\} | Y_1 \right] \right\} \right].$$

Finally, using our vector notation for paths and changing $e_1 \rightarrow e$ yields the expression in the claim. \square

2.4.1 The n -Stage Framework

We can easily generalize this framework to n stages, but we will need to write the expression recursively to keep it manageable. First, we assume we have a tuple of n RVs (Y_1, \dots, Y_n) that share a known joint-distribution $p_{WY_1 \dots Y_n}$ with W . For ease, define

$$Y_{\vec{i}} = (Y_1, \dots, Y_i),$$

$$\hat{W}_{\vec{i}} = \mathbb{E}[W|Y_{\vec{i}}].$$

Define the iterative sequence of functions $(J_i)_i$ by

$$J_1(Y_1) = \min_{e_i | \text{hd}(e_i)=s} \{E[J_2(Y_2, Y_1, e_1)|Y_1]\} \quad (2.8)$$

$$J_i(Y_i, Y_{i-1}, \{e_1, \dots, e_{i-1}\}) = \min_{e_i | \text{hd}(e_i)=\text{tl}(e_{i-1})} \{E[J_{i+1}(Y_{i+1}, Y_i, \{e_1, \dots, e_i\})|Y_i]\} \quad (2.9)$$

$$J_n(Y_n, Y_{n-1}, \{e_1, \dots, e_{n-1}\}) = \min_{p=(e_1, \dots, e_{n-1}, \dots)} \{p^T \hat{W}_{\bar{n}}\}. \quad (2.10)$$

The agent's average performance is

$$J(p_{WY_{\bar{n}}}) = E[J_1(Y_1)]. \quad (2.11)$$

2.4.2 Sequential Information Framework: Information Constraints

As in the non-sequential information case, we allow the agent to optimize the information it receives. For a tuple of capacities (C_1, \dots, C_n) , we restrict $p_{WY_{\bar{n}}}$ to a set $\Gamma(C_1, \dots, C_n)$. The agent's optimal performance under $\Gamma(C_1, \dots, C_n)$ is

$$J(C_1, \dots, C_n) = \min_{p_{WY_{\bar{n}}} \in \Gamma(C_1, \dots, C_n)} \{J(p_{WY_{\bar{n}}})\}. \quad (2.12)$$

We will define this set when we explore the sequential information case in more detail in Chapter 5.

2.5 Objectives

The objectives of this thesis are to:

1. provide computational and analytic upper and lower bounds for $J(C)$;
2. develop an algorithm that approximately or suboptimally solves $\min_{p_{WY} \in \Gamma(C)} \{J(p_{WY})\}$;
3. and efficiently compute analytic lower bounds for $J(C_1, \dots, C_n)$ by contrasting the performance of the sequential-information case to the non-sequential-information case.

We also present a generalization of our results to a multi-agent setting in Chapter 6, providing a limited study of the impact of information on network flow optimization.

2.6 Chapter Summary

We provided the notation, definitions, and basic concepts that will be used throughout this thesis. The specific frameworks for the value of sequential and non-sequential information in shortest path optimization were introduced and motivated by a general framework of partial information in stochastic optimization.

Chapter 3

Related Research

Several areas of study are related to our formulation and objectives, including the fields of order statistics, stochastic optimization, and shortest path optimization on random graphs. In this chapter, we provide a brief review of the existing works in these fields and detail their relationships to our framework.

3.1 Order Statistics and Optimization

3.1.1 Order Statistics

For a list of RVs $X = \{X_1, X_2, \dots, X_N\}$, we define the *minimum order statistic* as $X_{(1)} = \min_i \{X_i\}$. Under certain conditions, the PDF of $X_{(1)}$ can be computed using the identity presented in the following (well-known) proposition.

Proposition 5. *If $X_i \sim p_i$ and are independent, then $X_{(1)} \sim p$ with*

$$p(x) = \sum_i p_i(x) \prod_{j \neq i} P(X_j \geq x) \quad (3.1)$$

Proof. The event $X_{(1)} > x$ is equivalent to $\{X_i > x \text{ for all } i\}$. The PDF is computed by taking a derivative. \square

$E[X_{(1)}]$ can be computed directly from this identity when the X_i 's are independent. The computation is even simpler in the case of i.i.d. variables. However, our problem of computing $J(p_{\hat{W}})$ is a special case of computing $E[X_{(1)}]$ when the X_i 's are not generally independent:

$$J(p_{\hat{W}}) = E[X_{(1)}] = E\left[\min_p \{X_p\}\right]$$

where $X_p = \sum_{e \in p} \hat{W}_e$. In the non-i.i.d. case, computing $E[X_{(1)}]$ is difficult, and so bounds are typically computed instead. In this section, we review some methods for computing

these bounds and how they apply to computing bounds for $J(p_{\hat{W}})$.

3.1.2 Analytic Bounds based in Order Statistics

Bounds for $E[X_{(1)}]$ in the non-i.i.d. case are typically computed using only certain properties of the distribution, such as the means and variances of the individual RVs. One example of such a bound is the one presented in [1] that is computed by treating the min operator as a norm (by converting it to a max operator) and then bounding it by the 2-norm of the vector X . This method uses only first and second moment information about X . However, the bound is not directly amenable to sums of RVs.

An alternative type of bound is presented in [8] using Jensen's Inequality in combination with the moment generating functions for the X_i 's, but the method requires that the X_i 's be identically distributed, which does not apply to our framework.

A technique for bounding $E[X_{(1)}]$ that can be directly generalized to our case is presented in [5]. Though the bound itself is actually a low-complexity optimization, a conservative analytic expression can be derived simply by choosing a convenient suboptimal solution to the optimization.

We now present a bound for $J(C)$ based in a generalization of Theorem 1 in [5].

Corollary 1. $J(C) \geq J(0) - \frac{1}{2}\sqrt{|E|}C$.

Proof. First, we have

$$\begin{aligned} \min_p \{p^T \hat{W}\} &= \min_p \{p^T (\hat{W} - z) + p^T z\} \geq \min_p \left\{ \sum_{e \in p} \hat{W}_e - z_e \right\} + \min_p \{p^T z\} \\ &\geq \min_p \left\{ \sum_{e \in p} [\hat{W}_e - z_e]^- \right\} + \min_p \{p^T z\} \geq \sum_{e \in p} [\hat{W}_e - z_e]^- + \min_p \{p^T z\} \end{aligned}$$

The bound above is true for all z . Applying Proposition 1 and Theorem 2 in [5] to the above expression, we get

$$J(p_{\hat{W}}) \geq \max_{z_e} \left\{ \min_p \sum_{e \in p} z_e - \sum_e \left(E[\hat{W}_e] - z_e + \sqrt{(E[\hat{W}_e] - z_e)^2 + \text{VAR}[\hat{W}_e]} \right) \right\}.$$

Setting $z = \mu$ and applying Jensen's Inequality yields the lower bound in the claim. Since the bound is true for each $p_{\hat{W}}$, it is true for the minimizing such distribution and, thus, $J(C)$. Note that $J(0) = \min_p \left\{ \sum_{e \in p} \mu_e \right\}$. \square

While Corollary 1 provides an analytic lower bound for $J(C)$, the bound is independent of the graph's topology. Moreover, we will find that it can be easily derived from the analytic bounds in Chapter 4, which does include topology information and lends itself to information optimization.

3.1.3 Computational Lower Bounds via the Generalized Chebyshev's Inequality

Accurate bounds for non-i.i.d. order statistics can also be obtained using mathematical programming techniques. The basis for these approaches is a generalization of Chebyshev's Inequality, first presented in [15] and restated below.

Proposition 6 ([15]). *The following optimization over the distribution p_X*

$$\begin{aligned} \min_{p_X} E[f_0(X)] \quad & \text{subject to} \\ E[f_i(X)] &= z_i \quad \text{for } 1 \leq i \leq m \end{aligned}$$

is lower bounded by

$$\begin{aligned} \max_{y_0, \{y_i\}} \left\{ y_0 + \sum_i y_i z_i \right\} \quad & \text{subject to} \\ y_0 + \sum_{i=1}^m y_i f_i(x) &\leq f_0(x) \quad \text{for all } x \end{aligned}$$

Proof. Take the Lagrangian dual of the first optimization. □

Under certain conditions, there is no gap between the two optimizations in the proposition [15].

We can specialize Proposition 6 to computing a bound for $J(p_{\hat{W}})$ by setting

$$\begin{aligned} f_0(\hat{W}) &= \min_p \left\{ p^T \hat{W} \right\} \\ f_1(\hat{W}) &= \hat{W} \\ f_{e_1 e_2}(\hat{W}) &= \left(\hat{W} \hat{W}^T \right)_{e_1 e_2} \\ z_1 &= \mu \\ z_{e_1 e_2} &= \Lambda_{e_1 e_2} \end{aligned}$$

where Λ is the desired covariance matrix of \hat{W} . Applying Proposition 6, we get the following result.

Proposition 7.

$$\begin{aligned} J(p_{\hat{W}}) &\geq \max_{y_0, y, Y} \left\{ y_0 + \mu^T y + \text{Tr} \left(Y (\Lambda^2 + \mu \mu^T) \right) \right\} \quad \text{subject to} \\ \begin{bmatrix} Y & \frac{y - p}{2} \\ \frac{y^T - p^T}{2} & y_0 \end{bmatrix} &\geq 0 \quad \text{for all paths } p \end{aligned}$$

Proof. Applying Proposition 6 yields a quadratic constraint that can be rewritten as a semi-definite constraint [5]. □

We can compute a bound for $J(C)$ using a similar approach. We begin with the constraints

$$\begin{aligned} f_0(\hat{W}) &= \min_p \left\{ p^T \hat{W} \right\} \\ f_1(\hat{W}) &= \mathbb{E} \left[\hat{W} \right] \\ f_2(\hat{W}) &= \text{Tr} \left(\text{COV} \left[\hat{W} \right] \right) \\ z_1 &= \mu \\ z_2 &= C \end{aligned}$$

as well as an additional constraint $\text{VAR} \left[\hat{W}_e \right] \leq \sigma_e^2$. The inclusion of this inequality requires a slight change the form of the dual, but, for the purposes of this section, we do not expand on this treatment further.

Unfortunately, the dual optimizations for both $J(p_{\hat{W}})$ and $J(C)$ can have an exponential number of constraints (one for each path). A workaround is presented in [4], but to apply this technique, we need to relax our constraints f_i so that we constrain only the marginals distributions $\{p_{\hat{W}_e}\}$. We can apply this technique to the computation of $J(p_{\hat{W}})$ simply by setting

$$\begin{aligned} f_0(\hat{W}) &= \min_p \left\{ p^T \hat{W} \right\} \\ f_1(\hat{W}) &= \mathbb{E} \left[\hat{W} \right] \\ f_e(\hat{W}_e) &= \text{VAR} \left[\hat{W}_e \right] \\ z_1 &= \mu \\ z_e &= (\Lambda)_{ee} \end{aligned}$$

so that only the individual variances of \hat{W}_e are constrained. Of course, this allow the cross-covariances between the \hat{W}_e 's to be completely free.

Unfortunately, the capacity constraint $f_2(\hat{W}) = C$ in optimization for $J(C)$ disobeys the requirements on [4], and it cannot be relaxed without the resulting bound losing meaning. Nonetheless, the overall technique in [4] is useful, and we will present an extension of it in Chapter 4 as a comparison against our techniques.

Two drawbacks of the optimization-based approach taken in this section are the lack of analytic expressions for $J(C)$ and potential over-conservativeness. First, it is not immediately clear how one can derive analytic expressions for performance from the dual optimization, and so the impacts of graph topology and capacity on performance remain unclear using these techniques. Second, the expression for $J(C)$ only depends on the first and second marginal moments of the distributions, and the dependencies between the edge weight estimates generated by the optimization may not even be realizable, so it may not be appropriate to use as the basis for an information optimization algorithm.

3.2 Stochastic Optimization

Another related area of research is stochastic optimization, and numerous papers have considered the problem of bounding the improvement one gets from information in stochastic programming:

- [14] computes the value of full information using by subdividing the domain of the uncertainty and applying numeric algorithms over each domain.
- [2] assumes that $h(x, W)$ is concave in W and then leverage the resulting concavity of $J(W)$ to derive a numeric bound, but the bound is worse than that obtained using Jensen's Inequality.
- [10] considers the impact of partial information where partial information is represented by a signal that offers information about the underlying distribution of the uncertainty, and it is assumed that there are a finite number of such distributions.
- Finally, [3] considers a similar representation of partial information, but it seeks to determine the worst-case performance of the optimization over the unknown distributions.

In all, the related areas of study in stochastic optimization tend not to be specific to any particular application, and thus the bounds are computational, occasionally overly-conservative, and not amenable to an analytic analysis. Because we are considering a specific formulation (shortest path optimization), we seek to provide both computational and analytic bounds that are based in a common framework tied to the underlying structure of the problem.

3.3 Shortest Paths on Random Graphs

A final related area of study is that of shortest path algorithms on random graphs. However, existing bounds in this area tend to be either largely computational or they place heavy assumptions of the problem formulation to simplify analysis.

Two fairly general papers that can be used to generate bounds for $J(C)$ are [19] and [16]. However, the bounds are independent of the graph's topology. In fact, the bound in [16] is equal to a special case of our bound. Thus, neither bound is amenable to information optimization or the goals of this paper.

[18] and [13] compute bounds for performance that depend of topology and the edge weight distribution, but the assumptions in both are fairly heavy. [18] computes the PDF of the shortest path length in a complete graph having integer edge weights, and [13] studies of the average length of the shortest path in a complete graph with uniform edge weights. They do not generalize to our framework.

Finally, [9] computes a bound for arbitrary graphs having arbitrary distributions using dynamic programming that depends of the topology of the graph. The approach (in our case) produces a lower bound for $J(C)$. However, the algorithm is computational, and it is

also not clear how much information about the graph’s structure is actually utilized by the algorithm – by bounding performance along subpaths, some information about intersecting paths may be lost.

We will see that methods used in Chapter 4 supply bounds that depend on graph topology, are amenable to information optimization, and produce analytic bounds for performance.

3.4 Chapter Summary

We reviewed several related areas of study that share some overlap with either our problem formulation, our objectives, or the techniques used in this thesis. We found that none of the works mentioned in these areas address our objectives completely or are special cases of our results.

Chapter 4

The Value of Side Information in Shortest Path Optimization

In this chapter, we present new performance bounds for shortest path optimization under partial information that are amenable to both information optimization and analysis. These bounds are based on a new graph reduction that captures limited but significant information about the geometry of the graph’s path polytope. The resulting reduction will allow us to optimize information based on the portions of information that are relevant to determining the shortest path in the graph, not its length, and it will also allow us to capture the effects of the graph’s topology on performance.

4.1 Properties of the Path Polytope

Performance measures for graph algorithms are typically expressed in terms of limited characteristics of the graph. Such reductions are useful because they may capture a property of a graph that is either intuitive, easy to compute, or amenable for design, and yet have a significant effect on the algorithm’s performance. Examples of well-studied graph reductions include graph diameter, minimal cut, number and classifications of cliques, the spectrum of the matrix form of a graph’s matrix (adjacency or Laplacian), and so on.

A challenge in information optimization is the evaluation of the objective $J(p_{\hat{W}})$, which is known to be $\#P$ -hard [12]. To overcome this difficulty, we will concern our analysis with meaningful upper and lower bounds for $J(p_{\hat{W}})$ that can be efficiently computed. In this section, we present a graph reduction for shortest path optimization that will allow us to derive such bounds, and we will use these bounds to develop information optimization algorithms as well as analytic performance bounds.

4.1.1 Projection Matrix for the Path Polytope

The first property of \mathcal{P} that we explore is critical to the development of our low-complexity information optimization algorithms. First, let \bar{p} be the path of shortest average length,

defined as $\bar{p} = \operatorname{argmin}_{p \in \mathcal{P}} \{p^T \mu\}$. Clearly, $J(0) = \bar{p}^T \mu$.

Theorem 1. $\mathcal{P} - \bar{p}$ lies in a strict subspace $\mathcal{S}_{\mathcal{P}}$ of $\mathbb{R}^{|E|}$.

Proof. Because \mathcal{P} is a polytope, we only need to show that it does not have volume in $\mathbb{R}^{|E|}$.

Define $E_v = \{e \mid \operatorname{hd}(e) = v\}$ for $v \in V$. First, assume that $|E_s| > 1$. Define $H \in \mathbb{R}^{|E_s| \times |E|}$ by

$$H_{(i,j)} = \begin{cases} 1, & i = j \in E_s \\ 0, & \text{otherwise.} \end{cases}$$

Effectively, H is a diagonal projection matrix of $\mathbb{R}^{|E|}$ onto $\mathbb{R}^{|E_s|}$ in the sense that if we let $\hat{p} = Hp$, then $\hat{p}_e = p_e$ if $e \in E_s$ and is equal to zero otherwise.

By the virtue of G being DAG and s being the unique start vertex, any path p in G must contain exactly one of the edges in E_s . Hence, $HP = \{p \mid p_e = 1 \text{ for exactly one } e \in E_s \text{ and } p_e = 0 \text{ otherwise}\}$, and therefore HP is the simplex in $\mathbb{R}^{|E_s|}$.

The simplex (and, hence, HP) does not have volume in $\mathbb{R}^{|E_s|}$. Therefore, \mathcal{P} does not contain a hypercube \mathcal{C} of any size in $\mathbb{R}^{|E|}$ because if it did, $H\mathcal{C}$ would be a set with volume in the simplex. Because \mathcal{P} does not contain any hypercubes, it has no volume.

Now, assume that $|E_s| \leq 1$. If its equals zero, the claim is obviously true. If it equals one, select the first vertex v “after” s with $|E_v| > 1$ and apply the proof above to E_v . If $|E_v| = 1$ for all v (the only remaining case), then one can easily see that there is only one path in G , the vector $[1 \ 1 \ \dots \ 1]$, which is a single point in $\mathbb{R}^{|E|}$ and, thus, has no volume. \square

Because G is DAG, we can efficiently characterize $\mathcal{S}_{\mathcal{P}}$ by its projection matrix $H_{\mathcal{P}}$, and we can compute $H_{\mathcal{P}}$ in polynomial time. The details of the computation are not critical to the developments of this paper, so we save the details for the appendix. The following theorem formally states our claim.

Theorem 2. $H_{\mathcal{P}}$ can be computed in polynomial time.

4.1.2 Outer Spheric Approximation

The remaining portion of our graph reduction concerns the geometry of the boundary for \mathcal{P} . We simplify our description of the boundary by relying on low-complexity outer approximations, specifically a sphere. In general, finding the minimal-radius sphere containing a polytope is computationally hard [6], but in the case of the path polytope, it can be computed quite efficiently.

Theorem 3. The minimal-radius sphere $B(r_o^*, c_o^*)$ containing \mathcal{P} is given by the solution to

the convex quadratic optimization

$$\begin{aligned}
& \min r^2 \text{ subject to} \\
& r^2 \geq J(s) + \|c\|^2 \\
& \bar{J}(v) \geq \max_{e \mid \text{hd}(e)=v} \{\bar{J}(\text{tl}(e)) + (1 - 2c_e)\} \\
& \bar{J}(t) = 0.
\end{aligned} \tag{4.1}$$

Proof. Because the extreme points of \mathcal{P} are the paths p of the graph, a necessary and sufficient condition for a sphere with radius r and center c to contain \mathcal{P} is

$$\|p - c\|^2 \leq r^2 \Leftrightarrow p^T p - 2p^T c + c^T c \leq r^2.$$

However, since p is a 0-1 vector, $p^T p = 1^T p$, so we have the equivalent inequality

$$\begin{aligned}
& r^2 \geq p^T (1 - 2c) + \|c\|^2 \text{ for all } p \in P \\
& \Leftrightarrow r^2 \geq \max_{p \in P} \{p^T (1 - 2c)\} + \|c\|^2 \\
& \Leftrightarrow r^2 \geq \max_{p \in \mathcal{P}} \{p^T (1 - 2c)\} + \|c\|^2.
\end{aligned}$$

The expression $\max_{p \in \mathcal{P}} \{p^T (1 - 2c)\}$ equates to computing the longest path in an acyclic graph when the edge weights are given by $(1 - 2c)$. It is given by $\bar{J}(s)$ in the claim. \square

Corollary 2. $\mathcal{P} - c_o^* \subset \mathcal{S}_{\mathcal{P}}$.

Proof. The proof follows from the intuitively obvious (but not proven in this thesis) fact that $c_o^* \in \mathcal{P}$. \square

Remark 3. *Analytic bounds for (4.1) can be computed by first choosing a (possibly non-optimal) center c , computing the length of the longest path using under the edge weight vector $(1 - 2c)$, and then computing r . We will use this strategy in some later examples where we compute analytic performance bounds under a capacity constraint for certain interesting graph topologies.*

4.1.3 Inner Spheric Approximation

Efficient algorithms for generating inner spheric (as well as ellipsoidal) approximations to a polytope are well-known [6], and we do not reproduce these results in this paper. We do note, however, that such algorithms assume that the set to be approximated has volume, which \mathcal{P} does not. A way around this problem is to simply compute the inner approximation strictly within the affine subspace containing \mathcal{P} . The details of the computation are not important to the developments of this paper since we will not leverage inner spheric approximations beyond some basic results that are nearly identical to the case of outer spheric approximations.

Similar to the case of an outer spheric approximation, we let r_i^* and c_i^* respectively be the radius and center of the maximal inner sphere contained in \mathcal{P} .

4.2 Information Optimization using Graph Reductions

4.2.1 Information Optimization via Upper and Lower Bound Optimization

In this section, we apply the graph reductions of the last section to derive computational bounds for performance as well as information optimization algorithms. Our first bound follows from a direct application of inner and outer spheric approximations.

Lemma 1.

$$(c_o^*)^T \mu - r_o^* \max_{p_{\hat{W}} \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \hat{W}\| \right] \right\} \leq J(C) \leq (c_i^*)^T \mu - r_i^* \min_{p_{\hat{W}} \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \hat{W}\| \right] \right\}.$$

Proof. We start with the lower bound. Since $c_o^* \in \mathcal{P}$, $H_{\mathcal{P}}(\mathcal{P} - c_o^*) = \mathcal{P} - c_o^*$. Therefore,

$$\begin{aligned} J(p_{\hat{W}}) &\geq \mathbb{E} \left[\min_{p \in \mathcal{P} - c_o^*} \left\{ (p + c_o^*)^T \hat{W} \right\} \right] \\ &= \mathbb{E} \left[\min_{p \in H_{\mathcal{P}}(\mathcal{P} - c_o^*)} \left\{ p^T \hat{W} \right\} \right] + (c_o^*)^T \mu \\ &= \mathbb{E} \left[\min_{p \in (\mathcal{P} - c_o^*)} \left\{ (H_{\mathcal{P}} p)^T \hat{W} \right\} \right] + (c_o^*)^T \mu \\ &= \mathbb{E} \left[\min_{p \in (\mathcal{P} - c_o^*)} \left\{ p^T H_{\mathcal{P}} \hat{W} \right\} \right] + (c_o^*)^T \mu \\ &\geq \mathbb{E} \left[\min_{p \in B(r_o^*, 0)} \left\{ p^T H_{\mathcal{P}} \hat{W} \right\} \right] + (c_o^*)^T \mu \\ &= -r_o^* \mathbb{E} \left[\|H_{\mathcal{P}}(\hat{W})\| \right] + (c_o^*)^T \mu. \end{aligned}$$

The upper bound similarly follows. \square

Because evaluating $J(p_{\hat{W}})$ is impractical, we can consider optimizing information by optimizing these bounds. Specifically, we seek to maximize $\mathbb{E} \left[\|H_{\mathcal{P}} \hat{W}\| \right]$ over $\Gamma(C)$. The optimization can be performed using stochastic gradient descent.

Remark 4. An interpretation of Lemma 1 is that we should optimize side information by concentrating the energy of the estimate to the subspace $\mathcal{S}_{\mathcal{P}}$. The component of \hat{W} normal to that subspace (specifically, $(I - H_{\mathcal{P}})\hat{W}$) is lost in the projection. What does this component represent? It is the amount of length in \hat{W} common to all paths. Thus, it only aids in estimating the actual length of the paths. The tangential component $H_{\mathcal{P}}\hat{W}$ contains all of the information for path selection.

4.2.2 Upper Bound Optimization in the Gaussian Case

In the Gaussian case, upper bound $J(C)$ by a convex optimization. We begin with a useful background result.

Lemma 2. *Let $Z \sim N(0, I)$. The optimization*

$$\max_{\Lambda_{\hat{W}} \in \Gamma_G(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}$$

is equivalent to the convex optimization

$$\max_{\Lambda_{\hat{W}} \in \mathcal{L}} \left\{ \mathbb{E} \left[\|\sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}$$

where $\mathcal{L} = \Gamma_G(C) \cap \{\Lambda_{\hat{W}} \mid (H_{\mathcal{P}} - I)\Lambda_{\hat{W}} = 0\}$.

Proof. We have

$$\max_{\Lambda_{\hat{W}} \in \Gamma_G(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\} = \max_{\Lambda_{\hat{W}} \in \Gamma_G(C), \sqrt{\Lambda_{\hat{W}}} = H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}}} \left\{ \mathbb{E} \left[\|\sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}$$

where the equality comes from the fact that $H_{\mathcal{P}}H_{\mathcal{P}} = H_{\mathcal{P}}$.

Since $\text{Range}(\Lambda_{\hat{W}}) = \text{Range}(\sqrt{\Lambda_{\hat{W}}})$, we have the following equivalences for the constraint $\sqrt{\Lambda_{\hat{W}}} = H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}}$:

$$\begin{aligned} \sqrt{\Lambda_{\hat{W}}} = H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} &\Leftrightarrow \text{Range}(\sqrt{\Lambda_{\hat{W}}}) = \mathcal{S}_{\mathcal{P}} \Leftrightarrow \text{Range}(\Lambda_{\hat{W}}) = \mathcal{S}_{\mathcal{P}} \Leftrightarrow \Lambda_{\hat{W}} = H_{\mathcal{P}} \Lambda_{\hat{W}} \\ &\Leftrightarrow (H_{\mathcal{P}} - I)\Lambda_{\hat{W}} = 0. \end{aligned}$$

Clearly, this constraint is convex over $\Lambda_{\hat{W}}$. Therefore, \mathcal{L} is convex.

Now we prove the objective is concave. The function $f(\Lambda_{\hat{W}}) = z^T \Lambda_{\hat{W}} z = \|\sqrt{\Lambda_{\hat{W}}} z\|^2$ is linear over $\Lambda_{\hat{W}} \geq 0$. Therefore, $\sqrt{f(\Lambda_{\hat{W}})} = \|\sqrt{\Lambda_{\hat{W}}} z\|$ is concave for $\Lambda_{\hat{W}} \geq 0$. By linearity of the expected value operator, the objective is concave. \square

We now state a main result of this chapter.

Theorem 4. *If $p^T \mu \leq K$ for some K over all paths p , $J_G(C)$ is upper bounded by the convex optimization*

$$J_G(C) \leq K - r_i^* \max_{\Lambda_{\hat{W}} \in \mathcal{L}} \left\{ \mathbb{E} \left[\|\sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\} \quad (4.2)$$

where $Z \sim N(0, I)$, $\mathcal{L} = \Gamma_G(C) \cap \{\Lambda_{\hat{W}} \mid (H_{\mathcal{P}} - I)\Lambda_{\hat{W}} = 0\}$.

Proof. If $p^T \mu \leq K$, then

$$\begin{aligned} J_G(p_{\hat{W}}) &= \mathbb{E} \left[\min_p \left\{ p^T \hat{W} \right\} \right] \\ &\leq \mathbb{E} \left[\min_p \left\{ p^T \hat{W} + (K - p^T \mu) \right\} \right] \\ &= \mathbb{E} \left[\min_p \left\{ p^T (\hat{W} - \mu) \right\} \right] + K. \end{aligned}$$

By Corollary 1 and $\hat{W} \stackrel{d}{=} (\sqrt{\Lambda_{\hat{W}}} Z + \mu)$, we seek to compute

$$\max_{\Lambda \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}.$$

The remainder of the proof follows from Lemma 2. \square

4.2.3 Lower Bound Optimization in the Gaussian Case

We can also lower bound $J(C)$ by a convex optimization. We begin with some useful background results.

Proposition 8. *For any function $f : \text{paths} \rightarrow \mathbb{R}$ and any set $\overline{\mathcal{P}} \supset \mathcal{P}$,*

$$J(p_{\hat{W}}) \geq \min_{p \in \mathcal{P}} \{f(p)\} + \min_{p \in \overline{\mathcal{P}}} \left\{ p^T \hat{W} - f(p) \right\}.$$

Proof. The proof follows immediately from the fact that $\min_x \{a(x) + b(x)\} \geq \min_x \{a(x)\} + \min_x \{b(x)\}$. \square

Remark 5. *Proposition 8 provides a generalization of the approach taken in [5] to generate a low-complexity optimization for bounding the mean of the minimum order statistic. Let \mathcal{X} be the simplex in \mathbb{R}^n and let $W = (W_1, \dots, W_n)$ be a random vector in \mathbb{R}^n . Then the minimum order statistic is given by $\min_i \{W_i\} = \min_{x \in \mathcal{X}} \{x^T W\}$. We can derive the bound in [5] follows by setting $f(x) = x^T z$ for some vector $z \in \mathbb{R}^n$, setting $\overline{\mathcal{X}}$ to the unit cube, and then maximizing over z . A similar approach is taken in our previous work [17] for obtaining an analytic lower bound for $J(C)$.*

Corollary 3. $J(C) \geq J(0) - r_o^* \max_{p_{\hat{W}} \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}}(\hat{W} - \mu)\| \right] \right\}.$

Proof. Applying Proposition 8 with $f(p) = p^T \mu$ and $\overline{\mathcal{P}} = \mathcal{P}$ yields

$$\begin{aligned} J(p_{\hat{W}}) &\geq \min_{p \in \mathcal{P}} \{p^T \mu\} + \mathbb{E} \left[\min_{p \in \mathcal{P}} \left\{ p^T (\hat{W} - \mu) \right\} \right] \\ &= J(0) + \mathbb{E} \left[\min_{p \in \mathcal{P}} \left\{ p^T (\hat{W} - \mu) \right\} \right]. \end{aligned}$$

We can now apply Lemma 1 using $\hat{W} - \mu$ in place of \hat{W} yielding

$$J(C) \geq J(0) - r_o^* \max_{p_{\hat{W}} \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}}(\hat{W} - \mu)\| \right] \right\}.$$

□

Remark 6. *Note that both Proposition 8 and Corollary 3 are general results and are not specific to the Gaussian case.*

We now state the lower bound result.

Theorem 5. $J_G(C)$ is lower bounded by the convex optimization

$$J_G(C) \geq J(0) - r_o^* \max_{\Lambda \in \mathcal{L}} \left\{ \mathbb{E} \left[\|\sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}$$

where $Z \sim N(0, I)$, $\mathcal{L} = \Gamma_G(C) \cap \{\Lambda_{\hat{W}} \mid (H_{\mathcal{P}} - I)\Lambda_{\hat{W}} = 0\}$.

Proof. By Corollary 3,

$$J_G(C) \geq J(0) - r_o^* \max_{p_{WY} \in \Gamma(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}}(\hat{W} - \mu)\| \right] \right\}.$$

By Corollary 1 and our parameterization $\hat{W} = \sqrt{\Lambda_{\hat{W}}} Z + \mu$, we seek to compute

$$\max_{\Lambda_{\hat{W}} \in \Gamma_G(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\}.$$

The remainder of the proof follows from Lemma 2. □

4.2.4 Special Case Analytic Solution for Information Optimization in the Gaussian Case

Theorems 4 and 5 provide convex optimizations for information optimization. Under certain conditions, we can derive analytic expressions for the optimal performance as well as the corresponding optimizing covariance matrix $\Lambda_{\hat{W}}^*$.

Theorem 6. *Let $m = \dim(\mathcal{S}_{\mathcal{P}})$. If C is sufficiently small and the feasible set for $\Lambda_{\hat{W}}$ is $\Gamma_G(C) = \{\Lambda_{\hat{W}} \mid 0 \leq \Lambda_{\hat{W}} \leq \Lambda_W \text{ and } \text{Tr}(\Lambda_{\hat{W}}) \leq C\}$, then*

$$\max_{\Lambda_{\hat{W}} \in \Gamma_G(C)} \left\{ \mathbb{E} \left[\|H_{\mathcal{P}} \sqrt{\Lambda_{\hat{W}}} Z\| \right] \right\} = \sqrt{\frac{2}{\pi}} \sqrt{C}$$

with the optimizing $\Lambda_{\hat{W}}^* = U \Sigma U^T$ where U is a unitary matrix of the form $U = [x_1 \cdots x_m y_{m+1} \cdots y_{|E|}]$ with $\{x_i\}$ forming an orthonormal basis for $\mathcal{S}_{\mathcal{P}}$, and Σ is a diagonal matrix with $\Sigma_{ii} = C$ for some i and $\Sigma_{jj} = 0$ for $j \neq i$.

Proof. First, we can decompose $\Lambda_{\hat{W}}Z = Z^\parallel + Z^\perp$ where $H_{\mathcal{P}}Z^\parallel = Z^\parallel$ and $H_{\mathcal{P}}Z^\perp = 0$. By orthogonality, $\|\Lambda_{\hat{W}}Z\|^2 = \|Z^\parallel\|^2 + \|Z^\perp\|^2$. Therefore,

$$\|H_{\mathcal{P}}\Lambda_{\hat{W}}Z\| = \|Z^\parallel\| = \sqrt{\|\Lambda_{\hat{W}}Z\|^2 - \|Z^\perp\|^2} \leq \|\Lambda_{\hat{W}}Z\|.$$

Hence, $\Lambda_{\hat{W}}^*$ must satisfy $H_{\mathcal{P}}\Lambda_{\hat{W}}^*Z = \Lambda_{\hat{W}}^*Z$ in order to maximize the objective. We can parameterize the set of all feasible $\Lambda_{\hat{W}} \in \Gamma_G(C)$ satisfying this constraint as follows. Set $\Lambda_{\hat{W}} = U\Sigma U^T$ where the unitary matrix U has the form

$$U = [x_1 \cdots x_m y_{m+1} \cdots y_{|E|}]$$

with $\{x_i\}$ forming an orthonormal basis for $\mathcal{S}_{\mathcal{P}}$, and Σ is a diagonal matrix with $\Sigma_{ii} = C_i$ for $i \leq m$, $\Sigma_{ii} = 0$ for $i > m$, and nonnegative C_i satisfying $\sum_i C_i = C$. Optimizing $\Lambda_{\hat{W}}$ is now equivalent to optimizing over the basis $\{x_i\}$ and constants $\{C_i\}$.

Substituting this parameterization for $\Lambda_{\hat{W}}$ yields the optimization

$$\max_{\{x_i\}, \{C_i\}} \left\{ \mathbb{E} \left[\left\| \sum_{i=1}^m C_i x_i x_i^T Z \right\| \right] \right\} = \max_{\{x_i\}, \{C_i\}} \left\{ \mathbb{E} \left[\sqrt{\sum_{i=1}^m C_i (x_i^T Z)^2} \right] \right\}$$

where equality follows from the fact that the $\{x_i\}$ are an orthonormal basis.

By the concavity of $\sqrt{\cdot}$ and the fact that $\left\{\frac{C_i}{C}\right\}$ is on a simplex, we have

$$\begin{aligned} &= \max_{\{x_i\}, \{C_i\}} \left\{ \sqrt{C} \mathbb{E} \left[\sqrt{\sum_{i=1}^m \left(\frac{C_i}{C}\right) (x_i^T Z)^2} \right] \right\} \leq \max_{\{x_i\}, \{C_i\}} \left\{ \sqrt{C} \sum_{i=1}^m \left(\frac{C_i}{C}\right) \mathbb{E} \left[\sqrt{(x_i^T Z)^2} \right] \right\} \\ &= \max_{\{x_i\}, \{C_i\}} \left\{ \sqrt{C} \sum_{i=1}^m \left(\frac{C_i}{C}\right) K \right\} = \max_{\{x_i\}, \{C_i\}} \left\{ \sqrt{C} K \right\} = \sqrt{C} K \end{aligned}$$

where $K = \mathbb{E} \left[\sqrt{(n^T Z)^2} \right]$ for any unit vector n (by symmetry, any unit vector n yields the same number, and hence it is the same for each x_i). The upper bound is achieved if we choose any orthonormal basis $\{x_i\}$ and set $C_i = C$ for some i and $C_j = 0$ for all $j \neq i$. Note that the index i needs to be chosen so that $\Lambda_{\hat{W}} \leq \Lambda_W$. We assume C is sufficiently small in the theorem statement so that there exists such a i that satisfies this condition.

Finally, if we choose $n = [1 \ 0 \ \cdots \ 0]^T$, we get $K = \mathbb{E} |Z_1|$ where $|Z_1|$ has the folded normal distribution. Hence $K = \sqrt{\frac{2}{\pi}}$. \square

Corollary 4. *Under the conditions of Theorem 6, an optimal $\Lambda_{\hat{W}}^*$ is $\frac{C}{\|p_1 - p_2\|^2} (p_1 - p_2)(p_1 - p_2)^T$ for any two paths $p_1, p_2 \in P$.*

Proof. Using the definitions of Theorem 6, for the basis element x_i associated with the non-zero element of Σ , make $x_i = \frac{p_1 - p_2}{\|p_1 - p_2\|}$ since then $x_i \in \mathcal{S}_{\mathcal{P}}$. \square

Remark 7. Note that if we were to add additional restrictions to $\Gamma_G(C)$ such as $\Lambda_{\hat{W}} \sim \text{diagonal}$, it may not be practical to derive an analytic solution for the optimizing $\Lambda_{\hat{W}}^*$, and so we would be left with executing a convex optimization to solve for it.

Remark 8. An interesting property of the optimal $\Lambda_{\hat{W}}^*$ in Theorem 6 is that it concentrates the information to a single singular vector rather than spreading it over many singular vectors. Corollary 4 goes further and tells us that we should simply compare any two paths p_1 and p_2 of the graph over the edges for they do not intersect (the $p_1 - p_2$ expression). This is in qualitative agreement with Examples 1 and 2 where it was shown that concentrating information energy along a single path of a graph is superior to spreading it among many paths.

4.3 An Analytic Relationship between Capacity and Performance

4.3.1 Performance Lower Bounds

We can manipulate Corollary 3 to provide analytic lower bounds for performance in terms of the capacity C .

Theorem 7. $J(C) \geq J(0) - r_o^* \sqrt{C}$.

Proof. By $\sigma_{\max}(H_{\mathcal{P}}) = 1$ and Jensen's Inequality,

$$\mathbb{E} [\|H_{\mathcal{P}} \hat{W}\|] \leq \sqrt{\mathbb{E} [\|H_{\mathcal{P}} \hat{W}\|^2]} \leq \sqrt{\mathbb{E} [\|\hat{W}\|^2]} \leq \sqrt{C}$$

for all distributions $p_{\hat{W}} \in \Gamma(C)$. Simply apply this inequality to Corollary 3. \square

A lesson that can be drawn from Theorem 7 is that if the designer of the graph has some intuition about the relationship between the graph topology and the radius of the minimal sphere, then the designer should seek to maximize the radius. However, if the radius of the minimal sphere is not known and only a rough estimate of the benefit of information is desired, then the following topology-free bound can be applied.

Corollary 5. A proportionally-tight lower bound for $J(C)$ over all graph topologies is $J(C) \geq J(0) - \frac{1}{2} \sqrt{|E|} \sqrt{C}$.

Proof. Applying $c = [\frac{1}{2} \dots \frac{1}{2}]^T$ in the optimization in Theorem 3 sets $r = \frac{1}{2} \sqrt{|E|}$. Tightness is proven in Example 3. \square

Remark 9. The bound in Corollary 5 appears in [17] and [16] using different methods. In [17], it is obtained by bounding \mathcal{P} with the unit cube. In [16], it is obtained using convex majorization of RVs.

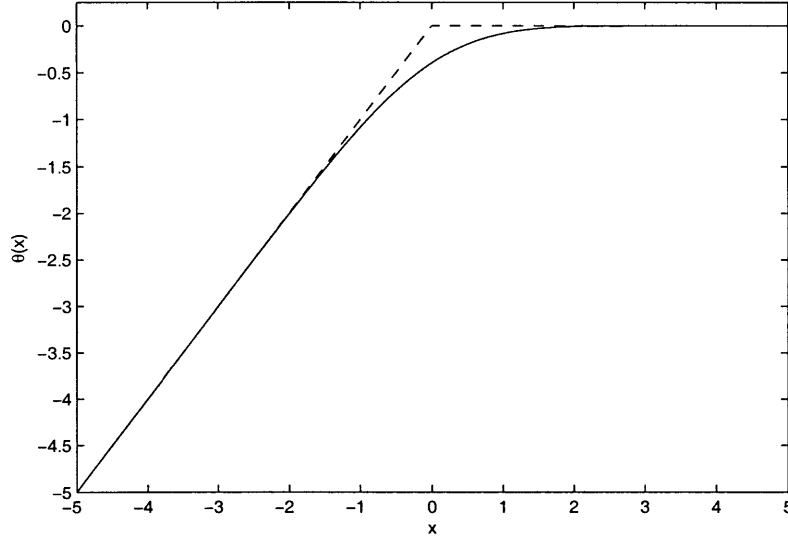


Figure 4-1: $\Theta(x)$ (solid line) compared to $\min\{x, 0\}$ (dashed line).

4.3.2 Performance Upper Bounds

An analytic upper bound for performance is more difficult to provide because the upper bound expressions we have so far derived do not seem to be expressible purely in terms of the capacity C (we cannot apply Jensen's Inequality as we did in the proof of Theorem 7). Moreover, restricted merely to the conditions $\mathbb{E}[\hat{W}] = \mu$ and $\text{VAR}[\hat{W}] \leq C$, a tight upper performance bound over distributions for any graph topology is $J(0)$.

Proposition 9. *A tight upper bound for $J(C)$ over all distributions subject to only a mean and variance constraint is $J(0)$ (the zero-information case).*

Proof. Without loss of generality, suppose $\mathbb{E}[\hat{W}] = 0$ and $\text{VAR}[\hat{W}] = 1$, and let

$$p_e^n(x) = \frac{1}{2n^2}\delta(x+n) + \frac{1}{2n^2}\delta(x-n) + \frac{2n^2-2}{2n^2}\delta(x).$$

Each distribution p_e^n satisfies the moment constraints, but it is easy to show that as $n \rightarrow \infty$, the limiting performance is $J(0)$. \square

To this end, we apply a different technique for producing performance upper bounds in terms of C for Gaussian edge weight estimates. The technique is based in reducing \mathcal{P} to a line.

First, let $\Theta_X(c) = \mathbb{E}[\min\{X, c\}]$. Clearly, $\Theta_X(c) \leq c$. Essentially, $\Theta_X(c)$ seeks to determine how much X is less than c on average. When $X \sim N(0, 1)$, we will simply write Θ without the subscript. In this case, $\Theta(0) = \frac{-1}{\sqrt{2\pi}}$.

Theorem 8. Assume independent edge weights and restrict $\Lambda \in \Gamma_G(C) \cap \{\Lambda \text{ is diagonal}\}$. For any $\hat{p} \in \mathcal{P}$, an upper bound for $J_G(C)$ is

$$J_G(C) \leq J(0) + \sqrt{\tilde{C}} \Theta \left(\frac{(\hat{p} - \bar{p})^T \mu}{\sqrt{\tilde{C}}} \right)$$

where $\tilde{C} = \min \{C, \text{VAR} [(\hat{p} - \bar{p})^T W]\}$.

Proof. Under independence assumptions, it is easy to verify that

$$\text{VAR} [(\hat{p} - \bar{p})^T \hat{W}] \leq \min \{C, \text{VAR} [(\hat{p} - \bar{p})^T W]\} = \tilde{C}.$$

It is also easy to construct a distribution $p_{\hat{W}}$ that corresponds to a diagonal $\Lambda \in \Gamma_G(C)$ that achieves this bound. Therefore, we have

$$(\hat{p} - \bar{p})^T \hat{W} \stackrel{d}{=} \sqrt{\tilde{C}} Z + (\hat{p} - \bar{p})^T \mu.$$

where $Z \sim N(0, 1)$.

Now,

$$\begin{aligned} J_G(\Lambda) &= \mathbb{E} \left[\min_{p \in \mathcal{P}} \{p^T \hat{W}\} \right] \leq \mathbb{E} \left[\min_{p \in \{\bar{p}, \hat{p}\}} \{p^T \hat{W}\} \right] \\ &= \mathbb{E} \left[\min \{ \bar{p}^T \hat{W}, \hat{p}^T \hat{W} \} \right] \\ &= \mathbb{E} \left[\bar{p}^T \hat{W} \right] + \mathbb{E} \left[\min \{ 0, (\hat{p} - \bar{p})^T \hat{W} \} \right] \\ &= J(0) + \mathbb{E} \left[\min \{ (\hat{p} - \bar{p})^T \hat{W}, 0 \} \right] \\ &= J(0) + \mathbb{E} \left[\min \{ \sqrt{\tilde{C}} Z + (\hat{p} - \bar{p})^T \mu, 0 \} \right] \\ &= J(0) + \mathbb{E} \left[\sqrt{\tilde{C}} Z \right] + \mathbb{E} \left[\min \{ (\hat{p} - \bar{p})^T \mu, -\sqrt{\tilde{C}} Z \} \right] \\ &= J(0) + 0 + \sqrt{\tilde{C}} \mathbb{E} \left[\min \left\{ \frac{(\hat{p} - \bar{p})^T \mu}{\sqrt{\tilde{C}}}, Z \right\} \right]. \end{aligned}$$

□

Corollary 6. Assume the conditions of Theorem 8 and further assume that $\hat{p}^T \mu = \bar{p}^T \mu$. An upper bound for $J_G(C)$ is

$$J_G(C) \leq J(0) - \frac{1}{\sqrt{2\pi}} \sqrt{\tilde{C}}$$

where $\tilde{C} = \min \{C, \text{VAR} [(\hat{p} - \bar{p})^T W]\}$.

4.4 Examples

4.4.1 Analytic Examples

We begin with several analytic examples that bound the value of side information under a capacity constraint for certain graph topologies.

Example 5. *In this example, we compute a lower analytic bound for performance for a binary tree with L levels. We first choose a center c for the outer sphere approximation, the process for which corresponds to selecting a center point c_e for each edge e . A binary tree consists of levels of edges, an edge's level being defined as its distance from the root vertex. We select c according to the scheme $c_e = 2^{-l_e}$ where l_e is the level of edge e . Because level l has 2^l edges in it, we have 2^l edges e with $c_e = 2^{-l}$.*

We compute a bound for r_o^ using this center as follows.*

$$\begin{aligned} (r_o^*)^2 &\leq \max_p \{p^T(1 - 2c)\} + \|c\|^2 \\ &= \left(1 - 2 \sum_{l=1}^L \frac{1}{2^l}\right) + \left(\sum_l \sum_{e \in \text{layer } l} 2^l \frac{1}{2^l}\right) \\ &= \frac{1}{2^{L-1}} - 1 + L \\ &\leq L. \end{aligned}$$

Therefore, the average performance under capacity C is lower bounded by $J(C) \geq J(0) - \sqrt{L}\sqrt{C}$.

Example 6. *We now compute a lower analytic bound for performance for a complete graph where the start vertex s and terminating vertex t are chosen randomly. Of course, since a complete graph is undirected, a direct application of (4.1) is not practical since it requires the computation the longest path in an undirected graph with non-negative edge weights, however we can compute a suboptimal analytic lower bound fairly easily.*

For a selection of s and t in the complete graph, we divide the edges of the graph into classes of edges and apply the same center point to each edge within the same class. The classes are

- $\mathcal{A} = \{\text{edge connecting } s \text{ and } t\},$
- $\mathcal{B} = \{\text{edges connecting } s \text{ to vertices } v \neq t\},$
- $\mathcal{C} = \{\text{edges connecting } v \neq s \text{ to } t\},$
- and $\mathcal{D} = \{\text{edges connecting } v \neq s \text{ to } u \neq t\}.$

We can show that $|\mathcal{A}| = 1$, $|\mathcal{B}| = |\mathcal{C}| = |V| - 2$, and $|\mathcal{D}| = \frac{(|V|-2)(|V|-3)}{2}$.

We assign the center points as follows:

- $e \in \mathcal{A} \Rightarrow c_e = a$,
- $e \in \mathcal{B} \Rightarrow c_e = b$,
- $e \in \mathcal{C} \Rightarrow c_e = c$,
- and $e \in \mathcal{D} \Rightarrow c_e = d$.

We can show that as long $0 \leq \{b, c, d\} \leq 1$,

$$r^2 \leq \max \left\{ 1 - 2a, \max_{2 \leq i \leq |V|-1} \{i - 2(b+c) - (i-2)d\} \right\} + \left(a^2 + (|V|-2)(b^2 + c^2) + \frac{(|V|-2)(|V|-3)}{2} d^2 \right).$$

Note that $((|V|-1) - 2(b+c) - (|V|-3)d)$ is the length of the acyclic path having the most edges (that is, the path corresponding to the graph's diameter).

If we set $a = \frac{1}{2}$ and $0 \leq \{b, c, d\} \leq 1$, we get

$$r^2 \leq ((|V|-1) - 2(b+c) - (|V|-3)d) + \left(a^2 + (|V|-2)(b^2 + c^2) + \frac{(|V|-2)(|V|-3)}{2} d^2 \right).$$

Maximizing over b , c , and d , we get $b = c = d = \frac{1}{|V|-2}$, and

$$\begin{aligned} (r_o^*)^2 &\leq \frac{4|V|^2 - 13|V| + 4}{4(|V|-2)} \\ &\approx |V| \\ &\approx \sqrt{|E|} \end{aligned}$$

where the final approximate equality is specific to the complete graph. Therefore, the average performance is lower bounded by $J(C) \gtrsim J(0) - \sqrt[4]{|E|}\sqrt{C}$.

Example 7. Our final analytic example computes a lower bound for average performance improvement on an Erdos-Renyi (E-R) random graph parameterized by (n, p) . For this example, we assume that the agent is only interested in traveling between connected vertices of the graph, and that it knows the realization of the graph but not the edge weights.

One can view an E-R graph as a generalization of a complete graph (the complete graph is always achieved if we set $p = 1$). For some connected component \hat{G} of the realization G , assign a start vertex s and a terminating vertex t , use the same classes of edges \mathcal{A} , \mathcal{B} , \mathcal{C} , and \mathcal{D} as in the complete graph case, and use the same values a , b , c , and d for those classes applied to the edges in this connected component.

Suppose there is at least one path from s to t in \hat{G} that is not the edge connecting s to t , and let $N_{\hat{G}}$ be the number of edges in the longest such path. Then using a derivation similar

to that of the complete graph, we get a bound on the radius $r_{\hat{G}}$ for that component:

$$r_{\hat{G}} \leq \sqrt[4]{N_{\hat{G}}} \sqrt{C}.$$

Maximizing this bound (which will serve to lower the performance lower bound) over all components corresponds to maximizing $N_{\hat{G}}$ over \hat{G} , the maximum of which is the diameter D_G of the graph.

Therefore, if there is at least one component in G with more than one acyclic path between a pair of two vertices in the component, we have

$$r_o^* \leq \sqrt[4]{D_G} \sqrt{C}.$$

Finally, if $\log n > np \rightarrow \infty$, then, as $np \rightarrow \infty$, we are guaranteed to have a giant component in the graph, and $D_G \approx \frac{\log n}{\log np}$ [7], yielding

$$\begin{aligned} \mathbb{E}[J(C)] &\gtrsim \mathbb{E}[J(0)] - \mathbb{E} \left[\sqrt[4]{D_G} \right] \sqrt{C} \\ &\geq \mathbb{E}[J(0)] - \sqrt[4]{\mathbb{E}[D_G]} \sqrt{C} \\ &\approx \mathbb{E}[J(0)] - \sqrt[4]{\frac{\log n}{\log np}} \sqrt{C}. \end{aligned}$$

4.4.2 Comparative Examples

We now compare our approach for computing performance bounds for $J(C)$ against an optimization designed for the same purpose. The optimization is based in the work of [4] for computing lower bounds for $J(p_{\hat{W}})$. Our extension of [4] will be equivalent to applying an outer approximation $\bar{\Gamma}(C)$ to $\Gamma(C)$ where $\bar{\Gamma}(C)$ is the set of all distributions $p_{\hat{W}}$ satisfying (a) $\mathbb{E}[\hat{W}] = \mu$, (b) $\text{VAR}[\hat{W}_e] \leq \sigma_e^2$, and (c) $\text{VAR}[\hat{W}] \leq C$.

To start, consider the following optimization:

$$\begin{aligned} \underline{J}(C) &= \min_{p_{\hat{W}}} \left\{ \mathbb{E} \left[\min_p \left\{ p^T \hat{W} \right\} \right] \right\} \text{ subject to} \\ &\mathbb{E}[\hat{W}] = \mu, \text{ VAR}[\hat{W}] \leq C \\ &0 \leq \text{VAR}[\hat{W}_e] \leq \sigma_e^2. \end{aligned} \tag{4.3}$$

Clearly, $J(C) \geq \underline{J}(C)$. We can solve (4.3) by solving the dual optimization, and there are conditions under which there is no duality gap [15]. However, the dual will have a constraint for each path in the graph, and, hence, may be impractical to solve for even a moderately-sized graph. An efficient approach to tackling optimizations similar to (4.3) is presented in [4], but it requires that the primal only contains constraints on the individual edge weight moments, which our capacity constraint clearly disobeys. Nonetheless, we can

adapt the approach in [4] to handle our capacity constraint quite nicely. We present this extension as a corollary to the main result Corollary 3.2 in [4]. The proof is contained in the appendix.

Corollary 7 (to Corollary 3.2 in [4]).

$$\begin{aligned} \underline{J}(C) = \min_{\{H_e\}, \{\lambda_e\}} & \left\{ \sum_e \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix} \cdot H_e \right\} \text{ subject to} \\ & H_e \leq 0, \quad H_e \geq - \begin{bmatrix} \lambda_e^2 + \mu_e^2 & \mu_e \\ \mu_e & 1 \end{bmatrix} \\ & \sum_e \left(\begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix} \cdot H_e \right) v_e \in \mathcal{P} \\ & 0 \leq \lambda_e^2 \leq \sigma_e^2, \quad \sum_e \lambda_e^2 \leq C \end{aligned}$$

where v_e is the elementary basis vector with the e^{th} component equal to 1.

Remark 10. While Corollary 7 will yield a tighter performance lower bound than the analytic bound in Theorem 7 (since the Theorem also only uses first and second moment information), it has several drawbacks. First, as we vary C , we need to re-solve the optimization in order to compute the resulting change in $\underline{J}(C)$. Second, it offers no intuition as to how graph topology or capacity impacts performance. Finally, we cannot use it to optimize information in the case of Gaussian edge weights since it only applies limited first and second moment information; it may yield a covariance matrix $\Lambda_{\hat{W}}$ that is not realizable.

We now present two examples comparing Theorem 7 to Corollary 7. In both, we consider two cases for the mean μ : (a) $\mu = 0$ and (b) a random μ . Because our analytic bounds grow unbounded for increasing C (even for $C > \text{VAR}[W]$), we compare the bounds only over $0 \leq C \leq \min_e \{\text{VAR}[W_e]\}$. This will allow us to examine the impact of increasing capacity without the negative impact of saturating information from any one edge. In both examples, the solution to Corollary 7 is computed using the MATLAB Toolbox CVX [11].

Example 8. In this example, we consider a graph having two disjoint paths from s to t and consisting of the same number of edges with each edge e satisfying $\text{VAR}[W_e] \geq 10$. We compute performance lower bounds over $0 \leq C \leq 10$. Figures 4-2a and 4-2b show the performances of two approaches in each case. For the case $\mu = 0$, the performances are identical over capacity and graph size. In the case of a random (non-zero) μ , the bound of Corollary 7 has slightly better performance, but it has roughly the same rate of performance improvement.

Example 9. In this example, we consider random DAGs consisting of ten vertices and with each edge e satisfying $\text{VAR}[W_e] \geq 10$. We compute performance lower bounds over $0 \leq C \leq 10$. Figures 4-2a and 4-2b show the performances of two approaches in each case.

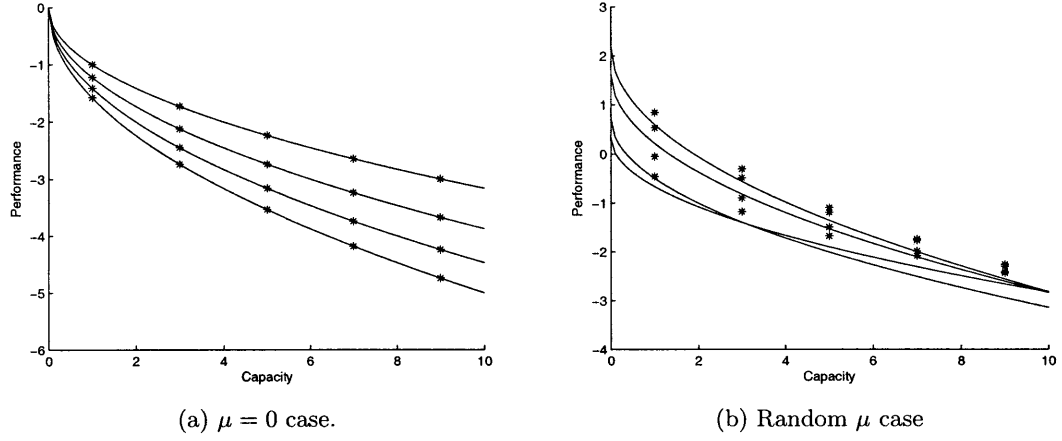


Figure 4-2: The analytic bound of Theorem 7 compared to the optimization-based bound of Corollary 7 for a two-path graph. The solid lines are the analytic bound performances, and the asterisks are the optimization-based bound performances. Each line represents a graph with an increasing number of links per path.

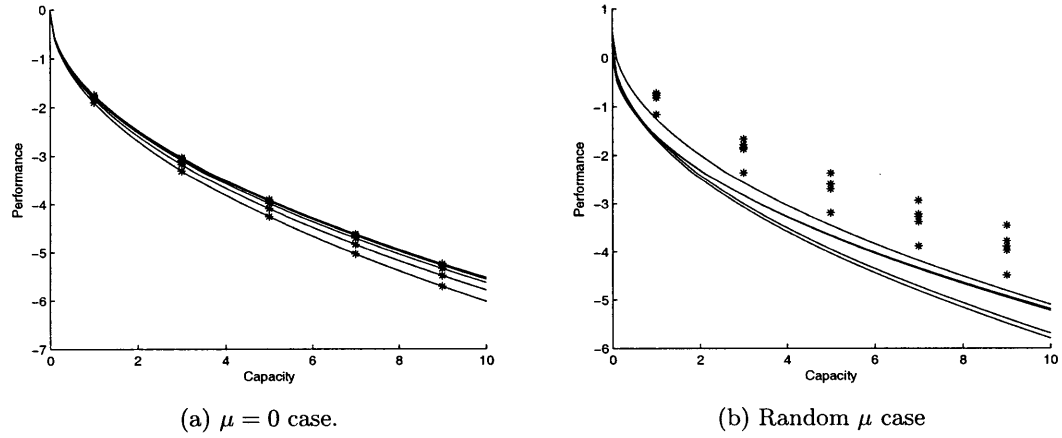


Figure 4-3: The analytic bound of Theorem 7 compared to the optimization-based bound of Corollary 7 for random DAGs. The solid lines are the analytic bound performances, and the asterisks are the optimization-based bound performances. Each line represents a different random graph topology with the number of vertices fixed.

For the case $\mu = 0$, the performances are identical over capacity and graph topology. In the case of a random (non-zero) μ , the bound of Corollary 7 has slightly better performance, but, once again, has roughly the same rate of performance improvement.

There are two interesting observations in the above examples. The first is that the bounds from both Theorem 7 and Corollary 7 seem to be identical in the regime $C \leq \min_e \{\text{VAR}[W_e]\}$ and $\mu = 0$, despite applying different outer-approximations of the constraint set. The second observation is that in the case of non-zero μ , the bounds at least seem to possess the same rate of improvement with increasing capacity.

4.5 Chapter Summary

A new graph reduction for analyzing the value of information for shortest path optimization was presented and used to provide algorithms for information optimization and analytic bounds for performance under a capacity constraint. It was shown that the graph reduction can be efficiently computed, and that information optimization in the case of a Gaussian distribution is bounded by a convex optimization. Examples comparing the analytic bounds of this chapter to an optimization-based bound for the same purpose showed similar (and, in some cases, the same) performance.

Chapter 5

The Value of Sequential Information in Shortest Path Optimization

In this chapter, we consider a generalization of the framework taken in the previous chapter by allowing the agent to receive information as it traverses the graph. We term this setup a *sequential-information* framework, and we consider two special cases of it: *controlled* and *uncontrolled* information. In the controlled information case, the agent can “optimize” the information it receives at each stage based on its location and past information. In the uncontrolled case, the information the agent receives at each stage is independent of the agent’s location and past information. The basis for our analysis is an abstraction of sequential decision-making.

5.1 Impact of Applying a Simple Information Constraint Set

Before we can proceed with an analysis of the sequential information cases, we need to define our information constraint set. First, we restate some notation that we first introduced in Chapter 2. For a tuple of capacities (C_1, \dots, C_n) , we restrict p_{WY_n} to a set $\Gamma(C_1, \dots, C_n)$. Define

$$\begin{aligned}\hat{W}_i &= \mathbb{E}[W|Y_i], \\ Y_{\vec{i}} &= (Y_1, \dots, Y_i), \\ \hat{W}_{\vec{i}} &= \mathbb{E}[W|Y_{\vec{i}}].\end{aligned}$$

For the $n = 2$ case, we define $\hat{W}_{12} = \mathbb{E}[W|Y_1Y_2]$ in order to make $\hat{W}_{\vec{2}} = \hat{W}_{12}$ more discernible \hat{W}_2 .

The agent's optimal performance under $\Gamma(C_1, \dots, C_n)$ is

$$J(C_1, \dots, C_n) = \min_{p_{WY_{\bar{n}}} \in \Gamma(C_1, \dots, C_n)} \{J(p_{WY_{\bar{n}}})\}. \quad (5.1)$$

Now, one may be tempted to consider the definition $\Gamma(C_1, C_2, \dots, C_n) = \Gamma(C_1 + C_2 + \dots + C_n)$, or, equivalently,

$$\left\{ p_{WY_{\bar{n}}} \mid \sum_{i=1}^n \text{VAR} [\hat{W}_i] \leq C \right\}$$

as our information constraint, but this definition has significant drawbacks. In particular, no matter how small (but positive) we make the C_i 's, the joint estimate $\hat{W}_{\bar{n}}$ may be arbitrarily accurate in the sense that $\text{VAR} [\hat{W}_{\bar{n}}] \approx \text{VAR} [W]$. The following proposition formalizes this claim.

Proposition 10. *For any $\epsilon > 0$,*

$$\begin{aligned} 0 = \inf_{p_{WY_{\bar{n}}}} \text{VAR} [W - \hat{W}_{\bar{n}}] \quad \text{subject to} \\ \text{VAR} [W] > 1 - \epsilon, \text{VAR} [\hat{W}_i] < \epsilon \end{aligned}$$

The same bound holds under the restriction $p_{Y_{\bar{n}}} = \prod_i p_{Y_i}$ (independence).

Proof. Without loss of generality, we prove the claim in the case of $n = 2$ and independent Y_1 and Y_2 . Define the RV W as a function of the RVs Y_1 and Y_2 by

$$W(Y_1, Y_2) = \begin{cases} h & (Y_1, Y_2) \in [0, \frac{1}{h}]^2 \\ 0 & \text{otherwise.} \end{cases}$$

It is straightforward to compute the following quantities:

$$\begin{aligned} \text{VAR} [W] &= 1 - \frac{1}{h^2} \\ \text{VAR} [\hat{W}_i] &= \frac{1}{h} - \frac{1}{h^2} < \frac{1}{h} \\ \text{VAR} [\hat{W}_{12}] &= \text{VAR} [W]. \end{aligned}$$

Setting h sufficiently large proves the claim. \square

An interpretation of Proposition 10 is that regardless of how little the “amount” of side information the agent receives at each time, the total side information can be nearly perfect. However, we are seeking the qualitative relationship “very little side information at each time gives very little total side information.” To this end, we will apply assumptions (or limitations) on the structure of the agent's information.

As we already mentioned, we consider two cases of sequential information: controlled and uncontrolled. In the controlled-information case, the agent can “optimize” the information it receives at each stage based on its location and past information. In the uncontrolled case, the information the agent receives at each stage is independent of the agent’s location and past information. We will develop different bounds for each of these cases that are partly based in the different definitions for $\Gamma(C_1, \dots, C_n)$ in each case.

5.2 The Uncontrolled Sequential Information Framework

We begin with an analysis of the uncontrolled information case. We will see that the resulting bound is a generalization of the performance bound in Chapter 4. In particular, if we set $C_1 = C$ and $C_i = 0$ for $i > 1$ (that is, asking for all the side information upfront), we get exactly the analytic lower bound in Chapter 4.

5.2.1 Structuring Information

We begin by defining our information sets $\{\Gamma(C_1, \dots, C_n)\}$. To motivate the general definition, we start with $n = 2$ and $\mu = 0$:

$$\Gamma(C_1, C_2) = \left\{ p_{WY_1Y_2} \mid \text{VAR} [\hat{W}_i] \leq C_i \text{ \& } \hat{W}_{12} = \hat{W}_1 + \hat{W}_2 \right\}.$$

This definition is reasonable in the sense that, in application, it is desirable to have a low-complexity estimation algorithm. Each time the agent receives information, computing the new estimate is simply a matter of adding the new estimate to the old one. This could correspond to, for example, the Y_i ’s having information about different edges weights. A consequence of this definition is that \hat{W}_1 and \hat{W}_2 must be independent.

In the general μ and n case, we want each partial estimation to have a similar structure as in the $n = 2$ case:

$$\Gamma(C_1, \dots, C_n) = \left\{ p_{WY_{\bar{n}}} \mid \text{VAR} [\hat{W}_i] \leq C_i \text{ \& } (\hat{W}_{\bar{k}} - \mu) = \sum_{i=1}^k (\hat{W}_i - \mu) \text{ for all } k \leq n \right\}.$$

The proposition below states a consequence of this information structure. We will not use the proposition directly to derive any further results, but it does serve to explicitly address our earlier concern about unbounded $\text{VAR} [W_{\bar{n}}]$.

Proposition 11. *For $p_{WY_{\bar{n}}} \in \Gamma(C_1, \dots, C_n)$, $\text{VAR} [\hat{W}_{\bar{n}}] \leq C_1 + \dots + C_n$.*

Proof. Without loss of generality, assume $\mu = 0$. Our constraint set Γ automatically yields the recursive relationship $\hat{W}_{\bar{k}} = \hat{W}_{\bar{k}-1} + \hat{W}_k$. Taking expectations given $Y_{\bar{k}-1}$ yields

$$\text{E} [W_{\bar{k}} | Y_{\bar{k}-1}] = \hat{W}_{\bar{k}-1} = \hat{W}_{\bar{k}-1} + \text{E} [\hat{W}_k | Y_{\bar{k}-1}].$$

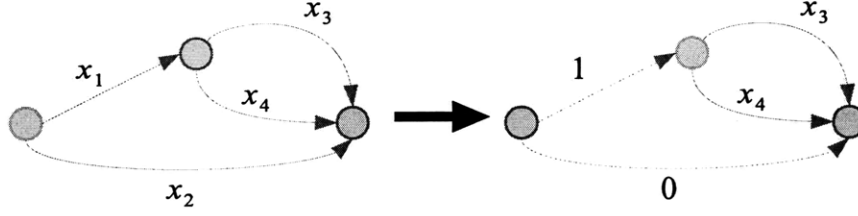


Figure 5-1: Example of a reduction in future decisions based on a past decision.

Hence $E[\hat{W}_k | Y_{k-1}] = 0$ and so \hat{W}_k is independent of Y_{k-1} , implying that it is also independent of $\hat{W}_{k-1} = \sum_{i=1}^{k-1} \hat{W}_i$. Therefore,

$$\text{VAR}[\hat{W}_k] \leq \text{VAR}[\hat{W}_{k-1}] + C_k.$$

Straightforward induction proves that $\text{VAR}[\hat{W}_k] \leq \sum_{i=1}^k C_i$. \square

5.2.2 The Geometry of Partial Decisions in the $n = 2$ Case

A challenge in computing $J(C_1, \dots, C_n)$ is the evaluation of $J(p_{WY_{\bar{n}}})$, which is known to be #P-hard [12], so, as in Chapter 4, we will instead focus on developing performance bounds for $J(p_{WY_{\bar{n}}})$ and $J(C_1, \dots, C_n)$. In this section, we develop the geometry underlying partial decisions that will be used to develop such bounds. We motivate our treatment with a simple example.

Example 10. Consider the left-side graph in Figure 5-1. Any path $p = [x_1 \ x_2 \ x_3 \ x_4]^T$ must be one of the vectors

$$p \in \{[1 \ 0 \ 1 \ 0]^T, [1 \ 0 \ 0 \ 1]^T, [0 \ 1 \ 0 \ 0]^T\} = P.$$

The agent starts at the leftmost vertex and can either choose the top or bottom path. In this example, it chooses the top path. From the next vertex, the set of paths that agent can choose from are

$$p \in \{[1 \ 0 \ 1 \ 0]^T, [1 \ 0 \ 0 \ 1]^T\} = P \cap \{[1 \ 0 \ x_3 \ x_4]^T \mid x_i \in \mathbb{R}\}.$$

The set $\{[1 \ 0 \ x_3 \ x_4]^T \mid x_i \in \mathbb{R}\}$ is an affine subspace in \mathbb{R}^4 .

Example 10 highlights the key qualitative concept in the geometry of partial decisions: a partial decision can be represented as an affine subspace S , the intersection of which with the decision space X is the set of remaining possible decisions: $S \cap X$.

We now apply this abstraction to our formulation in the two-stage case. We begin with several important definitions. Let

- $E_s = \{e \mid \text{hd}(e) = v\}$ be the set of edges connected to s ,

- $\hat{E}_e = \{\hat{e} \mid \text{there is no path } p \text{ with } p_e = p_{\hat{e}} = 1\}$ be the set of edges \hat{e} that cannot be reached from e ,
- $\bar{E}_e = \{\bar{e} \mid \text{there is no path } p \text{ with } p_e = 1, p_{\bar{e}} = 0\}$ be the set of edges \bar{e} that must be taken if e is taken,
- $S_e = \{x \in \mathbb{R}^{|E|} \mid x_e = 1, x_{\bar{e}} = 1 \ \forall \bar{e} \in \bar{E}_e, x_{\hat{e}} = 0 \ \forall \hat{e} \in \hat{E}_e\}$ be the affine subspace corresponding to choosing an edge $e \in E_s$,
- and $\mathcal{S} = \{S_e \mid x_e = 1 \text{ for exactly one } e \in E_s\}$ be the set of affine subspaces corresponding to the set of possible decisions at the start vertex s .

We can use our notation to express the two-stage performance in a simpler form.

Lemma 3.

$$J(p_{WY_1Y_2}) = \mathbb{E} \left[\min_{S \in \mathcal{S}} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \cap S} \left\{ p^T \hat{W}_{12} \right\} \mid Y_1 \right] \right\} \right]$$

Proof.

$$\begin{aligned} & \mathbb{E} \left[\min_{S \in \mathcal{S}} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \cap S} \left\{ p^T \hat{W}_{12} \right\} \mid Y_1 \right] \right\} \right] \\ &= \mathbb{E} \left[\min_{e \in E_s} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \cap S_e} \left\{ p^T \hat{W}_{12} \right\} \mid Y_1 \right] \right\} \right] \\ &= \mathbb{E} \left[\min_{e \in E_s} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \mid p_e = 1} \left\{ p^T \hat{W}_{12} \right\} \mid Y_1 \right] \right\} \right] \end{aligned}$$

□

5.2.3 Outer Approximating Partial Decisions in the $n = 2$ Case

A key step in our derivation of analytic performance bounds in the uncontrolled case is the use of a relaxation for partial decisions. The relaxation is an outer approximation of the set of affine subspaces corresponding to partial decisions. We begin by motivating the structure of the outer approximation for the two-stage case; that is, outer approximating \mathcal{S} .

Our outer approximation consists of larger set of affine subspaces that satisfy a looser set of constraints than \mathcal{S} while maintaining two important properties: (a) the dimension of the maximal-dimensional affine subspace $\max_{S \in \mathcal{S}} \{\dim S\}$, and (b) a particular geometric property of the intersection $S \cap B(0, r)$ for $S \in \mathcal{S}$.

We discuss the second property in more detail. A given $S \in \mathcal{S}$ has the form

$$S = \{x \in \mathbb{R}^{|E|} \mid x_e = 0 \text{ for } e \in E_0, x_e = 1 \text{ for } e \in E_1\}$$

where E_0 and E_1 are appropriately defined (but, of course, do not necessarily satisfy

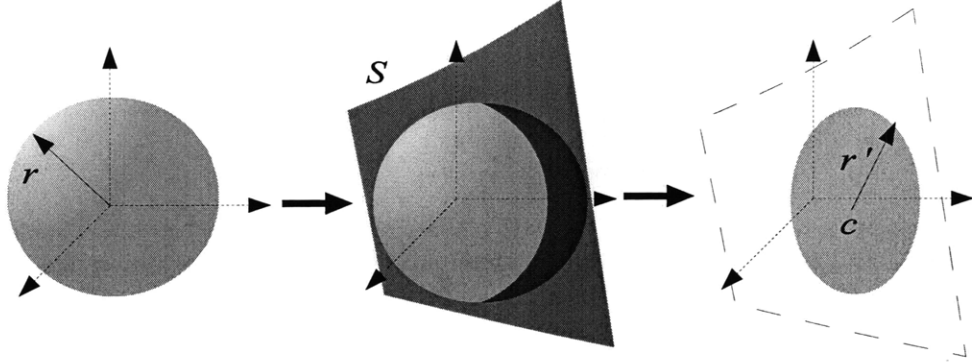


Figure 5-2: Illustration of a 3-dimensional ball with radius r intersecting an affine subspace S of dimension 2. The intersection is a 2-dimensional ball centered at c with radius r' (it appears as an ellipse due to distortion).

$E_0 \cup E_1 = E$). The intersection of S with the ball $B(0, r)$ yields a new set:

$$S \cap B(r, 0) = \left\{ x \in \mathbb{R}^{|E|} \mid x_e = 0 \text{ for } e \in E_0, x_e = 1 \text{ for } e \in E_1, \sum_{e \notin E_0, E_1} x_e^2 \leq \left(r^2 - \sum_{e \in E_1} 1^2 \right) \right\}.$$

From here, we can show that $S \cap B(0, r) = S \cap B(c, r')$ where

$$(c)_e = \begin{cases} 0, & e \notin E_0, E_1 \\ 0, & e \in E_0 \\ 1, & e \in E_1 \end{cases}$$

$$(r')^2 = r^2 - |E_1| = r^2 - \|c\|^2.$$

The key feature to notice is that $r' \leq r$ and $c \in S$. This leads us immediately to our outer approximation for \mathcal{S} .

$$\mathcal{S}^m = \{\text{Affine subspaces } S \mid \dim S = m, S \cap B(0, r) = S \cap B(c, r') \text{ for } r' \leq r, c \in S\}. \quad (5.2)$$

Figure 5-2 illustrates how an affine subspace $S \in \mathcal{S}^2$ operates on a ball in \mathbb{R}^3 .

Proposition 12. *If $\dim S \leq m$ for every $S \in \mathcal{S}$, then*

$$J(p_{WY_1Y_2}) = \mathbb{E} \left[\min_{S \in \mathcal{S}} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \cap S} \{p^T \hat{W}_{12}\} \mid Y_1 \right] \right\} \right] \geq \mathbb{E} \left[\min_{S \in \mathcal{S}^m} \left\{ \mathbb{E} \left[\min_{p \in \mathcal{P} \cap S} \{p^T \hat{W}_{12}\} \mid Y_1 \right] \right\} \right].$$

Proof. The proof follows from the fact that each $S \in \mathcal{S}$ is contained in some $S' \in \mathcal{S}^m$. \square

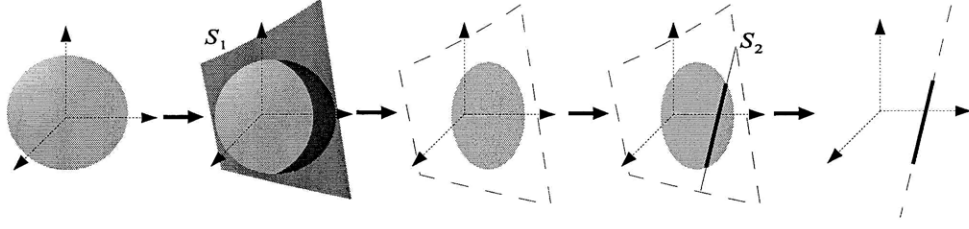


Figure 5-3: Illustration of a 3-dimensional ball intersecting an affine subspace S_1 of dimension 2 and then with another affine subspace $S_2 \subset S_1$ of dimension 1. With each intersection, we get a lower-dimensional ball (a line segment is a ball of dimension 1) with a new radius and center point.

5.2.4 Outer Approximating Partial Decisions in the General n Case

The general n -stage case relaxation follows immediately from the 2-stage case. First, define $S_1^{m_1} = S^{m_1}$ to be the set of affine subspaces for the first decision and define the set of affine subspaces for the i^{th} decision as

$$S_i^{m_i}(S_{i-1}) = \left\{ \text{Affine subspaces } S_i \mid \begin{array}{l} S_i \subset S_{i-1}, \dim S_i = m_i, S_i \cap B(0, r) = S_i \cap B(c_i, r') \text{ for } r' \leq r, c_i \in S_i \end{array} \right\} \quad (5.3)$$

where $S_j \in S_j^{m_j}(S_{j-1})$ for $j > 1$ and $S_1 \in S^{m_1}$, and the m_i 's are given. Each m_i represents the largest possible dimension of the remaining decision set at the i^{th} step. An important property in this general case is that the i^{th} affine subspace S_i must lie in the affine subspace S_{i-1} of the previous decision. Aside from this modification, the definition of the set is the same as in the two-stage case.

Figure 5-3 illustrates how a sequence of affine subspaces $S_1 \in S_1^2(\mathbb{R}^3)$ and $S_2 \in S_2^1(S_1)$ operate on a ball. Notice that each intersection yields a lower-dimensional ball with a new radius and center. Each successive intersection yields the set of remaining choices available to the agent at that step.

Once again, as an affine subspace, each $S_i \in S_i$ is completely characterized by (a) a projection matrix H_i (corresponding to its “centered” counterpart \hat{S}_i), and (b) an offset vector $c_i \in S_i$. We select c_i so that $\|c_i\| \leq \|x\|$ for all $x \in S_i$. By the Projection Theorem, c_i is orthogonal to \hat{S}_i . It is also intuitively clear that c_i is the center of the ball resulting from the intersection of S_i with a ball; that is, $S_i \cap B(0, r) = S_i \cap B(c_i, r')$. Hence, if we let B_i be the ball resulting from a sequence of intersections with a ball $B(0, r)$ in $\mathbb{R}^{|E|}$:

$$B_i = B(0, r) \cap S_1 \cap S_2 \cap \cdots \cap S_i = B_{i-1} \cap S_i$$

with $S_i \in S(S_{i-1})$, then c_i is the center of B_i . Of course, $B_0 = \mathbb{R}^{|E|}$.

Remark 11. Note that although c_i is in an affine subspace, the ambient space is still $\mathbb{R}^{|E|}$, and so we treat it as a vector in $\mathbb{R}^{|E|}$.

Of particular interest to us for analysis purposes is the difference $c_i - c_{i-1}$ in the center after each intersection. It will be convenient for us to separate this vector into its length and direction. Thus, we define

$$c_i = \frac{c_i - c_{i-1}}{\|c_i - c_{i-1}\|}$$

$$\gamma_i = \|c_i - c_{i-1}\|.$$

Finally, let

$$\gamma_i = \|c_i\|$$

be the distance from the origin to the center of ball B_i .

We now state a series of useful properties regarding these parameters. The first merely states that the center c_i of ball B_i can be written as the sum of the differences in the centers, which yields a convenient recursion for the centers.

Proposition 13. $c_i = \sum_{j=1}^i \gamma_j c_j = c_{i-1} + \gamma_i c_i$.

The second property states that c_i lies in the “centered” subspace \dot{S}_{i-1} but is orthogonal to the new “centered” subspace \dot{S}_i .

Proposition 14. $c_i \in \dot{S}_{i-1} \cap \dot{S}_i^\perp$.

Proof. We prove the claim for $c'_i = \gamma_i c_i = c_i - c_{i-1}$.

Since $S_i \subset S_{i-1}$, $\dot{S}_i \subset \dot{S}_{i-1}$, and, since c_{i-1} is orthogonal to \dot{S}_{i-1} , it is orthogonal to \dot{S}_i as well. By definition, c_i is orthogonal to \dot{S}_i . Therefore, since c'_i is a linear combination of the two, it is also orthogonal to \dot{S}_i .

Choose any $s \in S_{i-1}$. Since $c_{i-1} \in S_{i-1}$ and $c_i \in S_i \subset S_{i-1}$, we have $c_{i-1} = x + s$ and $c_i = y + s$ for some $x, y \in \dot{S}_{i-1}$. Thus, $c'_i = x - y \in \dot{S}_{i-1}$. \square

Now we prove an important relationship regarding c_i and the projection matrix H_i . The matrices H_i and $c_i c_i^T$ have orthogonal range spaces (since H_i is the projection matrix for \dot{S}_i and $c_i \perp \dot{S}_i$), but we can capture the range of the sum $H_i + c_i c_i^T$ in the previous subspace \dot{S}_{i-1} .

Proposition 15. $H_i + c_i c_i^T \leq H_{i-1}$.

Proof. To prove the claim, we show that $A = (H_i + c_i c_i^T)$ is a projection matrix for a subspace of \dot{S}_{i-1} since H_{i-1} is a projection matrix for \dot{S}_{i-1} .

Since c_i is a unit vector, $c_i c_i^T$ is a projection matrix. Because c_i is orthogonal to \dot{S}_i , $c_i^T H_i = [0 \cdots 0]$. It is easy to see that $A^T A = A$ and is thus a projection matrix.

Since H_i is a projection matrix for $\dot{S}_i \subset \dot{S}_{i-1}$, and since $c_i \in \dot{S}_{i-1}$, A is a projection matrix for a subspace of \dot{S}_{i-1} . \square

Finally, we have an important recursion relating γ_i to γ_{i-1} .

Proposition 16. $\gamma_i = \sqrt{\gamma_i^2 + \gamma_{i-1}^2} = \sqrt{\sum_{j=1}^i \gamma_j^2}$.

Proof. By definition,

$$\gamma_i^2 = \gamma_i^2 + \gamma_{i-1}^2 - 2c_i^T c_{i-1}.$$

Therefore, we only need to show that $c_i^T c_{i-1} = \gamma_{i-1}^2$. We can write $c_i = s + c_{i-1}$ for some $s \in \dot{S}_{i-1}$. By orthogonality, $s^T c_{i-1} = 0$, therefore $c_i^T c_{i-1} = \gamma_{i-1}^2$. \square

We will heavily leverage the equivalence between S_i and the parameterization (H_i, c_i) in our analysis. For a given sequence of subspaces, we can slightly modify this parameterization to

$$(S_1, S_2, \dots, S_i) = ((c_i, \gamma_1, H_1), (c_2, \gamma_2, H_2), \dots, (c_i, \gamma_i, H_i)).$$

5.2.5 Analytic Performance Lower Bound under Uncontrolled Information

We now present the main result for the uncontrolled information case – an analytic performance lower for the uncontrolled information case. Before stating the main result, we state a useful background result that allows us to remove the effect of μ from our bounds and apply our outer approximations.

Lemma 4. $J_1 \geq \underline{J}_1$ where

$$\begin{aligned} \underline{J}_1(Y_1) &= J(0) + \min_{S_1 \in \mathcal{S}_1} \{E[\underline{J}_2(Y_2, Y_1, S_1)|Y_1]\} \\ \underline{J}_i(Y_i, Y_{i-1}, S_{i-1}) &= \min_{S_i \in \mathcal{S}_i(S_{i-1})} \{E[\underline{J}_{i+1}(Y_{i+1}, Y_i, S_i)|Y_i]\} \\ \underline{J}_n(Y_n, Y_{n-1}, S_{n-1}) &= \min_{p \in \mathcal{P} \cap S_{n-1}} \left\{ p^T (\hat{W}_{\bar{n}} - \mu) \right\}. \end{aligned}$$

If $\mu = 0$, then $J_i \geq \underline{J}_i$.

Proof. Using the inequality $\min_x \{a(x) + b(x)\} \geq \min_x \{a(x)\} + \min_x \{b(x)\}$, we have

$$\begin{aligned} J_n((Y_n, Y_{n-1}, e_{n-1})) &\geq \min_{p \in P} \{p^T \mu\} + \min_{p \in P} \{p^T \hat{W}_{\bar{n}} - p^T \mu\} \\ &= J(0) + \min_{p \in P} \{p^T \hat{W}_{\bar{n}} - p^T \mu\}. \end{aligned}$$

It is clear that we can lower bound J_i by applying the outer approximation set \mathcal{S}_i . Thus,

we have $J_1 \geq \underline{J}'_1$ where

$$\begin{aligned}\underline{J}'_1(Y_1) &= \min_{S_1 \in \mathcal{S}_1} \{ \mathbb{E} [\underline{J}'_2(Y_2, Y_1, S_1) | Y_1] \} \\ \underline{J}'_i(Y_i, Y_{i-1}, S_{i-1}) &= \min_{S_i \in \mathcal{S}_i(S_{i-1})} \{ \mathbb{E} [\underline{J}'_{i+1}(Y_{i+1}, Y_i, S_i) | Y_i] \} \\ \underline{J}'_n(Y_n, Y_{n-1}, S_{n-1}) &= J(0) + \min_{p \in \overline{\mathcal{P}} \cap S_{n-1}} \left\{ p^T (\hat{W}_{\overline{n}} - \mu) \right\}.\end{aligned}$$

It is straightforward to verify that $\underline{J}_1 = \underline{J}'_1$ and, thus, $J_1 \geq \underline{J}_1$. If $\mu = 0$, it is clear that $J_i \geq \underline{J}_i$. \square

The main theorem for the uncontrolled information case is now given.

Theorem 9. *For a given distribution $p_{WY_{\overline{n}}} \in \Gamma(C_1, \dots, C_n)$ in the uncontrolled-information case, let Δ_i be the sum of the $|E| - m_{i-1}$ smallest singular values of $\text{COV} [\hat{W}_i]$. For any non-negative capacities $\sum_{i=1}^n C_i = C$,*

$$J(p_{WY_{\overline{n}}}) \geq J(0) - r_o^* \sqrt{C - \sum_{i=1}^n \Delta_i}.$$

Remark 12. *There is a straightforward interpretation of Theorem 9. $|E| - m_{i-1}$ represents the minimum number of edges that either cannot be traversed or must be traversed after the $i^{\text{th}} - 1$ stage; namely, those edges for which information at the i^{th} stage is not valuable because either (a) those edges must be taken anyway, or (b) those edges cannot be taken anymore. The effective capacity loss $\sum_{i=1}^n \Delta_i$ indicated in Theorem 9 represents the amount of information energy spread along these edges. If the information is broadcast to the agent or chosen randomly from the graph, an information spread to “useless” edges may be inevitable.*

Proof. The proof is by induction. We will use Proposition 22 in the appendix as the base case. By Lemma 4, we can assume without loss of generality that $\mu = 0$.

First, assume that for any $1 < i < n$,

$$\mathbb{E} [J_{i+1}(Y_{i+1}, Y_i, S_i) | Y_i] \geq -\sqrt{(r_o^*)^2 - \gamma_i^2} \left(\sum_{j=i+1}^n (C_j - \Delta_j) + \hat{W}_i^T H_i \hat{W}_i \right)^{\frac{1}{2}} + c_i^T \hat{W}_i.$$

We prove that the equivalent bound exists for J_i . To simplify notation, let

$$\begin{aligned}r^2 &= (r_o^*)^2 - \gamma_{i-1}^2, \\ K_1 &= \sum_{j=i+1}^n (C_j - \Delta_j) + \hat{W}_i^T H_i \hat{W}_i, \\ K_2 &= c_i^T \hat{W}_i.\end{aligned}$$

Using the relationship $\gamma_i^2 = \gamma_i^2 + \gamma_{i-1}^2$, we get

$$\mathbb{E} [J_{i+1}(Y_{i+1}, Y_i, S_i) \mid Y_i] \geq -\sqrt{r^2 + \gamma_i^2} \sqrt{K_1} + K_2.$$

Furthermore, since $c_i = \gamma_i c_i + c_{i-1}$, we can expand K_2 as

$$K_2 = \gamma_i c_i^T \hat{W}_i + c_{i-1}^T \hat{W}_i = \gamma_i c_i^T \hat{W}_i + K'_2,$$

where we select c_i so that $c_i^T \hat{W}_i \leq 0$ in order to further lower the bound.

By the lower bound in Lemma 4 for J_i , we now have the lower bound

$$\begin{aligned} & \mathbb{E} [J_i(Y_i, Y_{i-1}, S_{i-1}) \mid Y_{i-1}] \\ & \geq \mathbb{E} \left[\min_{(\gamma_i, c_i, H_i) \in \mathcal{S}_i(S_{i-1})} \left\{ -\sqrt{r^2 + \gamma_i^2} \sqrt{K_1} + \gamma_i c_i^T \hat{W}_i + K'_2 \right\} \mid Y_{i-1} \right]. \end{aligned}$$

Minimizing over γ_i (actually, over $\gamma_i \geq 0$ if we assume that $c_i^T \hat{W}_i \leq 0$) yields a new lower bound:

$$\mathbb{E} [J_i(Y_i, Y_{i-1}, S_{i-1}) \mid Y_{i-1}] \geq \mathbb{E} \left[\min_{(c_i, H_i) \in \mathcal{S}_i(S_{i-1})} \left\{ -\sqrt{r^2} \sqrt{K_1 + (c_i^T \hat{W}_i)^2} + K'_2 \right\} \mid Y_{i-1} \right].$$

We first examine the term $M = (c_i^T \hat{W}_i)^2$. Note that as M increases, we decrease the lower bound, so we seek an upper bound for M . Using the matrix inequality $c_i c_i^T \leq H_{i-1} - H_i$, we have

$$\begin{aligned} M &= \hat{W}_i c_i c_i^T \hat{W}_i \\ &\leq \hat{W}_i (H_{i-1} - H_i) \hat{W}_i \\ &= \left(\hat{W}_i H_{i-1} \hat{W}_i + 2 \hat{W}_i H_{i-1} \hat{W}_{i-1} + \hat{W}_{i-1} H_{i-1} \hat{W}_{i-1} \right) - \hat{W}_i H_i \hat{W}_i \end{aligned}$$

Substituting this upper bound for M yields

$$\begin{aligned} K_1 + M &\leq \left(\sum_{j=i+1}^n (C_j - \Delta_j) + \hat{W}_i^T H_i \hat{W}_i \right) + \left(\hat{W}_i H_{i-1} \hat{W}_i + 2 \hat{W}_i H_{i-1} \hat{W}_{i-1} + \hat{W}_{i-1} H_{i-1} \hat{W}_{i-1} \right) \\ &\quad - \hat{W}_i H_i \hat{W}_i \\ &= \sum_{j=i+1}^n (C_j - \Delta_j) + \hat{W}_i H_{i-1} \hat{W}_i + 2 \hat{W}_i H_{i-1} \hat{W}_{i-1} + \hat{W}_{i-1} H_{i-1} \hat{W}_{i-1} \\ &= K'_1. \end{aligned}$$

Substituting our expressions for K'_1 and K'_2 yields a new lower bound:

$$\mathbb{E} [J_i(Y_i, Y_{i-1}, S_{i-1}) \mid Y_{i-1}] \geq \mathbb{E} \left[\min_{(\gamma_i, c_i, H_i) \in \mathcal{S}_i(S_{i-1})} \left\{ -\sqrt{r^2} \sqrt{K'_1} \right\} + K'_2 \mid Y_{i-1} \right].$$

Because K'_1 does not depend on S_i (that is, it is independent of the minimization), we have

$$\begin{aligned} \mathbb{E} [J_i(Y_i, Y_{i-1}, S_{i-1}) \mid Y_{i-1}] &\geq \mathbb{E} \left[-\sqrt{r^2} \sqrt{K'_1 + K'_2} \mid Y_{i-1} \right] \\ &\geq -\sqrt{r^2} \sqrt{\mathbb{E} [K'_1 \mid Y_{i-1}] + \mathbb{E} [K'_2 \mid Y_{i-1}]} \end{aligned}$$

where the inequality comes from Jensen's Inequality.

We now seek to compute the expectations. First,

$$\mathbb{E} [K'_2 \mid Y_{i-1}] = \mathbb{E} \left[c_{i-1}^T \hat{W}_i \mid Y_{i-1} \right] = c_{i-1}^T \hat{W}_{i-1}.$$

As far as $\mathbb{E} [K'_1 \mid Y_{i-1}]$, we compute an upper bound. Since H_i is a projection matrix for \dot{S}_i , $N_i = I - H_i$ is a projection matrix for \dot{S}_i^\perp . Let $\{n_j^i\}_j$ be an orthonormal basis for \dot{S}_i^\perp . There are $|E| - m_i$ such vectors, and we can write $N_i = \sum_j n_j^i (n_j^i)^T$. Now, we have

$$\begin{aligned} \mathbb{E} [K'_1 \mid Y_{i-1}] &= \sum_{j=i+1}^n (C_j - \Delta_j) + \mathbb{E} \left[\hat{W}_i^T H_{i-1} \hat{W}_i \mid Y_{i-1} \right] + \\ &\quad 2 \mathbb{E} \left[\hat{W}_i^T \mid Y_{i-1} \right] H_{i-1} \hat{W}_{i-1} + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1} \\ &\stackrel{(a)}{=} \sum_{j=i+1}^n (C_j - \Delta_j) + \mathbb{E} \left[\hat{W}_i^T (I - N_{i-1}) \hat{W}_i \right] + 0 + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1} \\ &\stackrel{(b)}{\leq} \sum_{j=i+1}^n (C_j - \Delta_j) + C_i - \sum_j (n_j^i)^T \mathbb{E} \left[\hat{W}_i \hat{W}_i^T \right] n_j^i + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1} \\ &\stackrel{(c)}{\leq} \sum_{j=i+1}^n (C_j - \Delta_j) + C_i - \Delta_i + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1} \\ &= \sum_{j=i}^n (C_j - \Delta_j) + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1}. \end{aligned}$$

where (a) follows from the independence of \hat{W}_i and Y_{i-1} as well as the equality $H = I - N_{i-1}$, (b) follows from $\text{VAR} [\hat{W}_i] = C_i$ and our expression for N_{i-1} in terms of the vectors n_j^i , and (c) follows from the fact that the n_j^i 's are orthonormal vectors.

Therefore,

$$\begin{aligned} &\mathbb{E} [J_i(Y_i, Y_{i-1}, S_{i-1}) \mid Y_{i-1}] \\ &\geq -\sqrt{(r_o^*)^2 - \gamma_{i-1}^2} \left(\sum_{j=i}^n (C_j - \Delta_j) + \hat{W}_{i-1}^T H_{i-1} \hat{W}_{i-1} \right)^{\frac{1}{2}} + c_{i-1}^T \hat{W}_{i-1}. \end{aligned}$$

Now we prove the base case. Proposition 22 (see the appendix) provides the lower bound

$$\begin{aligned} \mathbb{E} [J_n(Y_n, Y_{n-1}, S_{n-1}) | Y_{n-1}] &\geq -\sqrt{(r_o^*)^2 - \gamma_{n-1}^2} \left(\mathbb{E} [\hat{W}_n^T H_{n-1} \hat{W}_n | Y_{n-1}] \right. \\ &\quad \left. + 2 \mathbb{E} [\hat{W}_n^T H_{n-1} \hat{W}_{n-1} | Y_{n-1}] + \hat{W}_{n-1}^T H_{n-1} \hat{W}_{n-1} \right)^{\frac{1}{2}} \\ &\quad + c_{n-1} \hat{W}_{n-1}. \end{aligned}$$

Using an argument similar to that above, we have

$$\mathbb{E} [J_n(Y_n, Y_{n-1}, S_{n-1}) | Y_{n-1}] \geq -\sqrt{(r_o^*)^2 - \gamma_{n-1}^2} \left(C_n - \Delta_n + \hat{W}_{n-1}^T H_{n-1} \hat{W}_{n-1} \right)^{\frac{1}{2}} + c_{n-1} \hat{W}_{n-1}.$$

To prove the claim, substitute $S_0 = \Re^{|E|} \Rightarrow H_0 = I \Rightarrow \Delta_1 = 0$, $\hat{W}_0 = \mathbb{E}[W] = 0$, and $\sum_{j=1}^n C_i = C$ into the expression for $\mathbb{E}[J_1(Y_1)]$. \square

5.3 The Controlled Sequential Information Framework

We now proceed with an analysis of the controlled information case. This case also provides a generalization of the non-sequential case but along different lines. In particular, to make any useful traction in this case, we need to apply a different information structure and performance analysis. For ease, we restrict our discussion in this section to the two-stage case, but it generalizes straightforwardly to the n -stage case. The key development in this section is not the bound itself, but rather practical conditions under which controlled requests may outperform non-sequential requests.

5.3.1 Structuring Information

As in the uncontrolled information case, we seek to structure the information the agent receives so that an accumulation of “little” bit of information at each time does not result in the agent knowing everything about the edge weights. However, we cannot apply the same information restriction used in the uncontrolled case to this case because, in that case, \hat{W}_1 and \hat{W}_2 were independent. In the controlled case, the agent’s location is used to determine future information, but its location is dependent on past information. Thus, there are dependencies among the estimates.

Intuitively, one may expect that the ability for the agent to concentrate future information to the relevant remaining portion of the graph should possibly result in improved performance relative to requesting all information upfront. However, in general, this intuition is not correct. The next proposition shows that if we were to simply constrain information so that $\text{VAR} [\hat{W}_{\bar{n}}] \leq C_1 + \dots + C_n$, then it is optimal to request all information upfront.

Proposition 17. *Let $Y = Y_{\bar{n}}$. For any nonnegative capacities C_i , $J(C_1, \dots, C_n) \geq J(C)$ where $C = \max_{p_{WY} \in \Gamma(C_1, \dots, C_n)} \left\{ \text{VAR} [\hat{W}(Y)] \right\}$.*

Proof. If $p_{WY_{\bar{n}}} \in \Gamma(C_1, \dots, C_n)$, then $J(p_{WY_{\bar{n}}}) \geq J(p_{WY})$ by an application of Jensen's Inequality on (2.7). Therefore, since C satisfies $\Gamma(C_1, \dots, C_n) \subset \Gamma(C)$,

$$\min_{p_{WY_{\bar{n}}} \in \Gamma(C_1, \dots, C_n)} \{J(p_{WY_{\bar{n}}})\} \geq \min_{p_{WY} \in \Gamma(C_1, \dots, C_n)} \{J(p_{WY})\} \geq \min_{p_{WY} \in \Gamma(C)} \{J(p_{WY})\}.$$

□

The basis for Proposition 17 is that under such a mild information restriction any joint distribution that can be achieved using controlled requests can also be achieved upfront. This implies that the type of information that the agent can request upfront can be very complicated and can even emulate sequential requests. Hence, we will seek to compare the performance of controlled requests to non-sequential requests in the case where the agent cannot emulate sequential requests upfront. The following example motivates the information structure that we will eventually apply to meet this criterion.

Example 11. Suppose $W \sim N(0, I)$ and that we restrict the agent to requesting whole edge weights. Therefore, we consider the case $\hat{W}_e \in \{W_e, 0\}$ and integer-valued C_i . We let the agent request C_1 edges weights on the first step, and C_2 different edge weights on the second step.

Let A_1 be a constant diagonal 0-1 matrix with at most C_1 number of 1's. We have $\hat{W}_1 = A_1 W$. Let A_2 be another constant diagonal 0-1 matrix with at most C_2 number of 1's. The new information the agent gets is $A_2 W$. Because the agent selects different edge weights at the second time step, A_1 and A_2 must have orthogonal range spaces or, equivalently, $A_1^T A_2 = 0$. Finally, because the selection of edge weights at the second step depends on information at the first step, we let A_2 be a function of \hat{W}_1 and write $A_2 = A_{2|\hat{W}_1} = A_{2|A_1 W}$. Thus, we can interpret A_2 as a RV.

We cannot write $\hat{W}_2 = A_2 W$ since the very realization of A_2 gives us information about \hat{W}_1 , and thus \hat{W}_2 actually has information outside the range space of A_2 . However, we know that our information after the second step is $A_1 W + A_2 W$. Thus, we have $\hat{W}_{12} = A_1 W + A_2 W$. A weaker condition that we satisfy is $\hat{W}_{12} \stackrel{d}{=} A_1 W + A_2 W$.

There are two key properties for us in the above example. The first is the orthogonal range space condition $A_1^T A_2 = 0$, and the second is the distribution of the joint estimate $\hat{W}_{12} \stackrel{d}{=} A_1 W + A_2 W$. We use these two conditions as the basis for the information structure we apply in the controlled information case.

$$\begin{aligned} \Gamma(C_1, C_2) = \{p_{WY_1 Y_2} \mid (\hat{W}_{12} - \mu) &\stackrel{d}{=} A_1 Z + A_{2|A_1 Z} Z \\ \text{VAR}[A_1 Z] &\leq C_1, \text{VAR}[A_{2|A_1 Z} Z] \leq C_2 \\ A_1^T A_2 &= 0\} \end{aligned}$$

where Z is an RV, A_1 is a constant matrix, and A_2 is a matrix function of $A_1 Z$. From here on, we simplify notation by treating A_2 as a random matrix related to the RV $A_1 Z$.

This information structure limits the types of distributions that agent can use to obtain information, and it guarantees that “little” amounts of information do not accumulate to “too much” information.

Proposition 18. *For $p_{WY_1Y_2} \in \Gamma(C_1, C_2)$, $\text{VAR} [\hat{W}_{12}] \leq C_1 + C_2$.*

Proof.

$$\begin{aligned} \text{VAR} [W_{12}] &= \text{VAR} [A_1Z + A_2Z] \\ &= \text{VAR} [A_1Z] + \text{VAR} [A_2Z] + 2 \text{E} [W^T A_1^T A_2Z] \\ &\leq C_1 + C_2 + 0. \end{aligned}$$

□

5.3.2 Performance Bounds

We now state the main theorem of this section – a comparison of the controlled information case to the non-sequential information case.

Theorem 10. *For $p_{WY_1Y_2} \in \Gamma(C_1, C_2)$, a lower bound for the controlled information case is*

$$J(p_{WY_1Y_2}) \geq J(0) - r_o^* \sqrt{C_c},$$

and for $p_{WY} \in \Gamma(C_1 + C_2, 0)$, a lower bound for the non-sequential information case is

$$J(p_{WY}) \geq J(0) - r_o^* \sqrt{C_{ns}},$$

where

$$\begin{aligned} C_c &= \max_{A_1} \left\{ \text{E} \left[\max_{A_2} \{ \text{VAR} [H_{\mathcal{P}}(A_1Z + A_2Z) \mid A_1Z] \} \right] \right\}, \\ C_{ns} &= \max_{A_1} \{ \text{VAR} [H_{\mathcal{P}} A_1Z] \}. \end{aligned}$$

Furthermore, if A_1^ is the optimizing matrix for C_{ns} , and if there is a projection matrix H satisfying $\text{VAR} [H A_1^* Z] = C_1$, then $C_c \geq C_{ns}$.*

Remark 13. *In Chapter 4, it is shown that only the component $H_{\mathcal{P}} \hat{W}$ contains useful information about the shortest path in the graph; in other words, $(I - H_{\mathcal{P}}) \hat{W}$ cannot be used to determine which path is the shortest. Therefore, we should apply all of the information energy into the subspace of the path polytope. One can interpret Theorem 10 as a statement that the agent’s performance in the controlled information case can be better than in the non-sequential case because the agent may be better able to concentrate the energy of its estimates to the subspace. This is reflected in the inequality $C_c \geq C_{ns}$.*

Proof. Without loss of generality, assume that $E[W] = 0$. Write the affine subspace containing \mathcal{P} as $S_{\mathcal{P}} = \dot{S}_{\mathcal{P}} + c_{\mathcal{P}}$. This gives us

$$\begin{aligned}
J(p_{WY_1Y_2}) &= E \left[\min_{S \in \mathcal{S}} \left\{ E \left[\min_{p \in \mathcal{P} \cap S} \left\{ p^T \hat{W}_{12} \right\} | Y_1 \right] \right\} \right] \\
&\stackrel{(a)}{=} E \left[\min_{S \in \mathcal{S}} \left\{ E \left[\min_{p \in \mathcal{P} \cap S} \left\{ (p - c_{\mathcal{P}})^T \hat{W}_{12} \right\} | Y_1 \right] \right\} \right] + E[c_{\mathcal{P}}^T W] \\
&\stackrel{(b)}{=} E \left[\min_{S \in \mathcal{S}} \left\{ E \left[\min_{p \in \mathcal{P} \cap S} \left\{ (p - c_{\mathcal{P}})^T H_{\mathcal{P}} \hat{W}_{12} \right\} | Y_1 \right] \right\} \right] + 0 \\
&\stackrel{(c)}{=} E \left[\min_{S \in \mathcal{S}} \left\{ E \left[\min_{p \in \mathcal{P} \cap S} \left\{ p^T H_{\mathcal{P}} \hat{W}_{12} \right\} | Y_1 \right] \right\} \right] \\
&\stackrel{(d)}{\geq} E \left[\min_{S \in \mathcal{S}} \left\{ \min_{p \in \mathcal{P} \cap S} \left\{ p^T H_{\mathcal{P}} \hat{W}_{12} \right\} \right\} \right] \\
&\stackrel{(e)}{=} E \left[\min_{p \in \mathcal{P}} \left\{ p^T H_{\mathcal{P}} (A_1 Z + A_2 Z) \right\} \right]
\end{aligned}$$

where (a) follows from the fact that $c_{\mathcal{P}}$ is a constant independent of p and e_1 , (b) follows from the fact that $(p - c_{\mathcal{P}}) \in \dot{S}_{\mathcal{P}}$ so that $(p - c_{\mathcal{P}}) = H_{\mathcal{P}}(p - c_{\mathcal{P}})$, (c) follows from the fact that $c_{\mathcal{P}}^T H_{\mathcal{P}} = (H_{\mathcal{P}} c_{\mathcal{P}})^T = 0^T$, (d) follows from Jensen's inequality, and (e) follows from substituting the equivalent distribution $\hat{W}_{12} \stackrel{d}{=} A_1 Z + A_2 Z$.

From Theorem 7, we know that $J(C) \geq J(0) - r_o^* \sqrt{\text{VAR}[\hat{W}]}$. In this case,

$$\hat{W} = H_{\mathcal{P}}(A_1 Z + A_2 Z).$$

Therefore, for a fixed selection of A_1 and A_2 , the agent's performance is

$$J(C_1, C_2) \geq J(0) - r_o^* \sqrt{\text{VAR}[H_{\mathcal{P}}(A_1 Z + A_2 Z)]}.$$

Now, in the controlled information case, the agent can optimize A_2 depending on the realization of $A_1 Z$, hence the definition of C_c . In the non-sequential case, the agent can only optimize A_1 .

Finally, we prove the final part of the claim. If we set $A_1 = H A_1^*$ and $A_2 = (I - H) A_1^*$, then

$$\begin{aligned}
\text{VAR}[A_1 Z + A_2 Z] &= \text{VAR}[H A_1^* Z] + \text{VAR}[(I - H) A_1^* Z] + 2 \text{COV}[H A_1^* Z, (I - H) A_1^* Z] \\
&\leq C_1 + (\text{VAR}[A_1^* Z] + \text{VAR}[H A_1^* Z] - 2 \text{COV}[A_1^* Z, H A_1^* Z]) + 0 \\
&\leq C_1 + (C_1 + C_2) + \text{VAR}[H A_1^* Z] - 2 \text{VAR}[H A_1^* Z] \\
&= C_1 + (C_1 + C_2) - C_1 \\
&= C_1 + C_2,
\end{aligned}$$

and $A_1 Z + A_2 Z = A_1^* Z$. Therefore, $C_c \geq C_{ns}$. \square



Figure 5-4: Graph topology where sequential information is worse than non-sequential information.

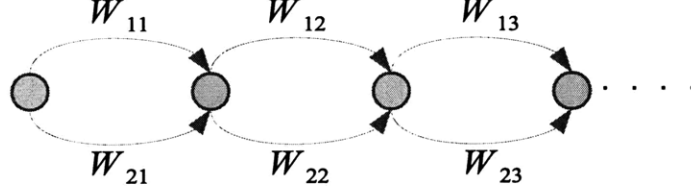


Figure 5-5: Graph topology where (controlled) sequential information can achieve the same performance as the non-sequential case.

5.3.3 Examples

We now present some examples that show how graph topology and our restriction on $\Gamma(C_1, C_2)$ impact whether controlled requests are beneficial relative to non-sequential requests.

Example 12. Consider the graph in Figure 5-4. The agent can take one of the two edges from the leftmost vertex. After traversing one of those edges, the remainder of its path is determined. Thus, the only useful information about W is contained in the first two edges, and so any information transmitted at the second stage is useless. Regardless of the restrictions on $\Gamma(C_1, C_2)$, $J(C_1, C_2) \geq J(C_1 + C_2)$.

Example 13. We reconsider Example 3 from Chapter 2. Consider the graph in Figure 5-5 with $W_{ij} \sim N(0, 1)$ and independent. Assume $C < \frac{|E|}{2}$, and suppose we restrict $p_{WY} \in \Gamma(C)$ so that $\hat{W} \sim N(\mu, \Lambda)$ with Λ diagonal (that is, we want our estimates of the edge weights to be independent Gaussians, just like W). In Example 3, it is shown that an optimal such p_{WY} satisfies (a) $\text{VAR}[\hat{W}_{1j}] = \frac{C}{|E|/2}$ and (b) $\text{VAR}[\hat{W}_{2j}] = 0$. Of course, by symmetry, one can immediately see that there are at least two such optimal distributions simply by reversing the roles of \hat{W}_{1j} and \hat{W}_{2j} .

It is clear that the agent only requires information about edge weights \hat{W}_{1j} and \hat{W}_{2j} immediately before having to decide which of the two to take (that is, at the j^{th} vertex). For the controlled information case, we take $p_{WY_1Y_2} \in \Gamma(\frac{C}{|E|/2}, C - \frac{C}{|E|/2})$ so that

1. $\text{VAR}[E[W_{11}|Y_1]] = \frac{C}{|E|/2}$ and $\text{VAR}[E[W_{ij}|Y_1]] = 0$ for $j > 1$,
2. and $\text{VAR}[E[W_{1j}|Y_1Y_2]] = \frac{C}{|E|/2}$ and $\text{VAR}[E[W_{2j}|Y_1Y_2]] = 0$ for $j > 1$.

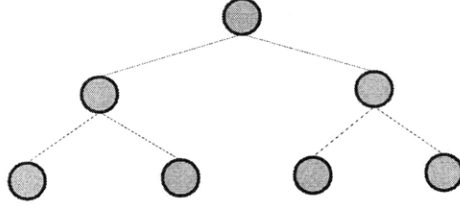


Figure 5-6: Graph topology where controlled sequential information can outperform non-sequential information.

In this case, it is straightforward to see that $J(p_{WY_1Y_2}) = J(C)$. In fact, we can show that $J\left(\frac{C}{|E|/2}, C - \frac{C}{|E|/2}\right) = J(C)$, meaning that $p_{WY_1Y_2}$ is optimal.

To map this setting to our controlled information framework, we set:

- $Z = W$,
- $A_1 \sim \text{diagonal with } (A_1)_{11} = \sqrt{\frac{C}{|E|/2}} \text{ and } (A_1)_{ii} = 0 \text{ for } i > 1$,
- and $A_2 \sim \text{diagonal with } (A_2)_{ii} = \sqrt{\frac{C}{|E|/2}} \text{ and } (A_2)_{11} = 0$.

One can easily check that $\hat{W}_{12} \stackrel{d}{=} A_1W + A_2W$.

Example 14. For our final example, we consider the graph in Figure 5-6. In this case, we assume $W \sim N(0, I)$ and $C_1 = C_2 = \frac{1}{2}$. We restrict $p_{WY_1Y_2} \in \Gamma(C_1, C_2)$ so that $\text{VAR}[E[W_e|Y_1]] \in \{0, \frac{1}{2}, 1\}$ and $\text{VAR}[E[W_e|Y_1Y_2]] \in \{0, \frac{1}{2}, 1\}$. In our setup, this is equivalent to setting $Z = W$ and restricting A_1 and A_2 to be diagonal matrices with $(A)_{ij} \in \{0, \sqrt{\frac{1}{2}}, 1\}$. Essentially, we want to restrict the agent to either getting perfect information about one edge or “half” information from each of two edges.

We now compute the agent’s performance under all of the possible request schemes for $p_{WY_1} \in \Gamma(C_1 + C_2, 0) = \Gamma(1, 0)$ (the non-sequential case):

Request Scheme	$J(p_{WY_1})$
Perfect information about any one edge	−0.3982
Partial information about two edges coming from the same vertex	−0.3983
Partial information about any two edges in series	−0.4808
Partial information about any two completely disconnected edges	−0.4819

Now, for the controlled information case, we consider only one information scheme. For ease, we simply describe the scheme:

1. First, request partial information about one of the two first edges coming from the root vertex. Call this edge e .

2. If $E[W_e|Y_1] \leq 0$, then traverse e , request partial information about one of the two remaining relevant edges, and take the one with the shortest estimated length.
3. If $E[W_e|Y_1] > 0$, then traverse the other edge, request partial information about one of the two remaining relevant edges, and take the one with the shortest estimated length.

Under this information scheme, we get $J(p_{WY_1Y_2}) \approx -0.5658$, which is better than all of the non-sequential information schemes.

5.4 Chapter Summary

In this chapter, we compared two types of sequential information schemes, controlled and uncontrolled, to the non-sequential information case. In each case, we introduced a practical assumption of the information's structure that led to meaningful comparisons of the cases. In the uncontrolled case, the lower bound for performance is higher than in the non-sequential case due to lost information energy that is spread across irrelevant edges of the graph. In the controlled case, if the information structure is sufficiently restricted, the lower bound for performance is better than in the non-sequential case.

Chapter 6

The Value of Information in Network Flow Optimization

We now study the impact of partial information in a social welfare context. We consider multiple agents that simultaneously choose decisions from a decision set and where the quality of each decision is dictated by (a) the decision itself, (b) some random perturbation, and (c) the number of agents that choose that same decision.

We motivate this problem with an example. Consider a traffic system where each vehicle corresponds to an agent. The decision set in this case is the set of available paths, and the quality of a decision is the total time of travel. However, the delay of a path is partially dictated by the number of agents who choose to take that path. Therefore, the agents must coordinate themselves in a non-greedy fashion if they wish to maximize their social welfare. We seek to determine how side information can benefit social welfare.

We begin our presentation with a specialization of the general framework in Chapter 2 to social welfare optimization. We then study this framework in the context of a specific problem: network flow optimization. For this application, we provide a specific relaxation that will allow us to express information optimization (on a lower performance bound) as a linear program as well as allow us to derive an analytic performance lower bound.

6.1 Partial Information in Stochastic Social Welfare Optimization

6.1.1 Performance

We first define a framework for studying the value of information in stochastic social welfare optimization.

Let $R \geq 1$ be the number of agents (an integer). Each agent makes a decision $x \in X$. For ease, we assume that X is finite. Let f_x be the fraction of the R agents who choose decision x . Define the vector of decisions as $f = [f_1 \dots f_{|X|}]^T$. Clearly, Rf_x is the number of agents who chose decision x , and Rf is the vector of how many agents chose each decision.

Finally, let $h(x, Rf, W)$ be the quality of decision x . We let the quality be a function of the collective decisions of the agents (the vector Rf) as well as some random perturbation W .

Denote the simplex in \mathbb{R}^n as $\Delta^n = \{x \in \mathbb{R}^n \mid x \geq 0, \|x\|_1 = 1\}$. Clearly, $f \in \Delta^{|X|}$, though by the assumption that R is an integer, f lies in a strict subset of the simplex.

Remark 14. *One can almost immediately see that this framework is actually a specialization of the stochastic optimization in Chapter 2. We have only extended our decision set to the set of all the agents' decisions.*

We assume that every agent has the same information Y about W . For a given Y , the agents optimize the average social performance

$$\min_f \mathbb{E} \left[\sum_{x \in X} f_x h(x, Rf, W) \mid Y \right] = \min_f \sum_{x \in X} f_x \mathbb{E} [h(x, Rf, W) \mid Y].$$

We denote the average social performance under the joint distribution p_{WY} as $J_R(p_{WY})$. This is simply given by

$$J_R(p_{WY}) = \mathbb{E} \left[\min_f \left\{ \sum_{x \in X} f_x \mathbb{E} [h(x, Rf, W) \mid Y] \right\} \right]. \quad (6.1)$$

6.1.2 Quantifying Information and Information Optimization

Finally, given a family of information constraint sets $\{\Gamma(C)\}$, the optimal performance under a capacity C is written as

$$J_R(C) = \min_{p_{WY} \in \Gamma(C)} \{J_R(p_{WY})\}. \quad (6.2)$$

6.2 Network Flow Optimization under Limited Information

In this section, we specialize the general stochastic social welfare framework to network flow optimization.

6.2.1 Specialization to Network Flow Optimization

In network flow optimization, we assume that the agents simultaneously traverse the graph, and that the edge weights of the graph are impacted by the number of agents on each edge as well as some inherent randomness.

The decision set is the set of paths P . Each agent chooses a path $p \in P$ to traverse. f_p is the fraction of agents who take path p , and, as a vector, $f \in \Delta^{|P|}$. We call f the *flow*. The quality $h(p, Rf, W)$ of path p is the total weight of that path, which is sum of

the weights of the edges along that path, and it is defined as

$$h(p, Rf, W) = \sum_{e \in p} (\alpha_e Rf'_e + W_e)$$

where

- $f'_e = \sum_{p \mid e \in p} f_p$ is the fraction of agents that specifically use edge e ,
- Rf'_e is the number of agents that use edge e ,
- $\alpha_e \geq 0$ is the sensitivity of the edge e 's weight to the number of agents using it,
- and W_e is an additive random weight for edge e .

We have some information Y about the edge weights W . The estimate of h given Y is

$$\mathbb{E}[h(p, Rf, W)|Y] = \sum_{e \in p} (\alpha_e Rf'_e + \hat{W}_e)$$

where $\hat{W} = \mathbb{E}[W|Y]$.

For a given Y , the agents choose the flow that optimizes the social welfare according to

$$\min_f \left\{ \sum_p f_p \mathbb{E}[h(p, Rf, W)|Y] \right\}.$$

The average performance of the agents under the joint distribution p_{WY} is now just

$$J_R(p_{\hat{W}}) = J_R(p_{WY}) = \mathbb{E} \left[\min_f \left\{ \sum_p f_p \sum_{e \in p} [\alpha_e Rf'_e + \hat{W}_e] \right\} \right]. \quad (6.3)$$

Note that we can denote our distribution p_{WY} as $p_{\hat{W}}$ since \hat{W} is the only random quantity in the expression.

6.2.2 Information Constraints

We define our information constraint sets $\{\Gamma(C)\}$ using the same definition in Chapter 2 for shortest path optimization:

$$\Gamma(C) = \left\{ p_{WY} \mid \text{VAR}[\hat{W}] \leq C \right\}. \quad (6.4)$$

6.3 Outer Approximating Network Flow Optimization

A key step in generating bounds for J_R is to relax the constraints on the flow f . A similar approach is taken in Chapter 4 for determining the value of information in shortest path

optimization. In Chapter 4, the constraint set for the agent is the set of paths P . By relaxing this constraint set, performance bounds could be derived. A significant difference between network flow optimization and shortest path optimization, however, is that shortest path optimization is a linear optimization while network flow optimization is a quadratic optimization. The methods used in Chapter 4 very much rely on its linear structure, so they are not useful to us here. To this end, we will develop a different approach that leverages the structure of the underlying formulation. Our method requires a minor assumption on the graph's structure.

Definition 3. *An edge cut across a graph is a collection of edges $\{e_1, e_2, \dots\}$ such that each path in the graph passes through exactly one edge in the collection.*

Assumption 3. *There is a subset $\mathcal{S} = \{S_i\}$ of all possible sets of edge cuts across G such that (a) each path p contains an edge from every cut $S_i \in \mathcal{S}$, and (b) for two different $S_i, S_j \in \mathcal{S}$, S_i and S_j are disjoint.*

For ease, let $N = |\{S_i\}|$. We expand the set of decisions for each agent using these sets.

Proposition 19. $P \subset \bar{P} = S_1 \times S_2 \times \dots \times S_N$.

Proof. By Assumption 3 and acyclicity, every path p has exactly one edge in each S_i . \square

Remark 15. *In the shortest path problem, the set of decisions for the agent is the set of paths P of the graph. A useful outer approximation for P is the unit hypercube in $\mathbb{R}^{|P|}$, which is the convex hull of all 1-0 combinations of edges. The hypercube approximation effectively expands the set of paths to all possible paths constructed from all the edges (including a path of zero edges that always has zero length). The outer approximation \bar{P} that we take for network flow optimization is similar – it is all combinations of edges subject to having only one edge from each cut.*

Proposition 20.

$$J_R(p_{\hat{W}}) \geq \sum_i \mathbb{E} \left[\min_{f \in \Delta_i} \left\{ \sum_{e \in S_i} f_e \left[\alpha_e R f_e + \hat{W}_e \right] \right\} \right]$$

where $\Delta_i = \Delta^{|S_i|}$ and $f \in \Delta_i$ is a vector $f = [f_{e_1} \ f_{e_2} \ \dots \ f_{e_{|S_i|}}]^T$ where $e_k \in S_i$.

Proof. The set of decisions \bar{P} corresponds to a graph, thus we can write the performance of the network flow formulation for this new graph as a lower bound for $J(p_{\hat{W}})$.

$$J_R(p_{\hat{W}}) \geq \mathbb{E} \left[\min_f \left\{ \sum_{p \in \bar{P}} f_p \sum_{e \in p} \left[\alpha_e R f'_e + \hat{W}_e \right] \right\} \right].$$

Clearly, $f \in \Delta^{|\bar{P}|}$. We can lower bound the optimization further if we allow f to take any value in this simplex (this is akin to making R a real number).

$$\geq \mathbb{E} \left[\min_{f \in \Delta^{|\bar{P}|}} \left\{ \sum_{p \in \bar{P}} f_p \sum_{e \in p} [\alpha_e R f'_e + \hat{W}_e] \right\} \right].$$

Now we reverse summations.

$$\begin{aligned} &\geq \mathbb{E} \left[\min_{f \in \Delta^{|\bar{P}|}} \left\{ \sum_e \left([\alpha_e R f'_e + \hat{W}_e] \sum_{p \mid e \in p} f_p \right) \right\} \right] \\ &= \mathbb{E} \left[\min_{f \in \Delta^{|\bar{P}|}} \left\{ \sum_e [\alpha_e R f'_e + \hat{W}_e] f'_e \right\} \right] \\ &= \mathbb{E} \left[\min_{f \in \Delta^{|\bar{P}|}} \left\{ \sum_i \sum_{e \in S_i} [\alpha_e R f'_e + \hat{W}_e] f'_e \right\} \right] \end{aligned}$$

We can equivalently optimize over f'_e directly, so we only need to determine the range of its value. Define the vector $f_i \in \mathcal{R}^{|S_i|}$ as $(f_i)_e = f'_e$ for $e \in S_i$. f_i is simply the vector of flow on edges in cut S_i . Because all flow must pass through edge cut S_i , it is clear that $f_i \in \Delta_i$, yielding

$$\begin{aligned} &= \mathbb{E} \left[\min_{\{f_i\} \in \Delta_i} \left\{ \sum_i \sum_{e \in S_i} [\alpha_e R (f_i)_e + \hat{W}_e] (f_i)_e \right\} \right] \\ &= \mathbb{E} \left[\sum_i \min_{f_i \in \Delta_i} \left\{ \sum_{e \in S_i} [\alpha_e R (f_i)_e + \hat{W}_e] (f_i)_e \right\} \right] \\ &= \mathbb{E} \left[\sum_i \min_{f \in \Delta_i} \left\{ \sum_{e \in S_i} [\alpha_e R f_e + \hat{W}_e] f_e \right\} \right] \end{aligned}$$

□

6.4 Information Optimization via Lower Bound Optimization

We now derive computational bounds for performance that are used to provide a rudimentary information optimization algorithm.

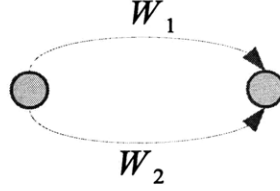


Figure 6-1: Graph for Example 15

6.4.1 The Impact of Information Selection: Examples

Before deriving our first set of bounds, we first provide some examples that highlight how information can impact performance.

Example 15. Consider the graph in Figure 6-1. Set $\alpha_1 = \alpha_2 = \alpha$ and let $W_i \sim N(0, 1)$ and independent. We consider $p_{\hat{W}} \in \Gamma(C)$ for $C \leq 2$ subject to the additional restriction that \hat{W}_1 and \hat{W}_2 are independent Gaussians as well. We write $\hat{W}_i \sim N(0, \lambda_i^2)$.

Evaluating (6.3) is a matter of solving for the optimal flows f_1 and f_2 corresponding to each of the two paths. It is straightforward, albeit lengthy, to compute $J(p_{\hat{W}})$ for this case, so we only provide the final expression for it:

$$J_R(p_{\hat{W}}) = J_R(0) + \mathbb{E} \left[\frac{C^2}{\alpha R} [Z]^{\frac{\alpha R}{C}} \left[[Z]^{\frac{\alpha R}{C}} - Z \right] \right]$$

where $Z = \frac{\hat{W}_2 - \hat{W}_1}{C}$ ¹. Unfortunately, even in this simple case, the average performance expression is far too difficult to exactly analyze. With some effort, though, we can lower bound as

$$J_R(C) \geq J_R(0) + 2C\phi\left(\frac{\alpha R}{C}\right) - \frac{C^2}{\alpha R},$$

where ϕ is the density of the normal distribution. Though simpler, this expression is not much of an improvement for the purposes of developing a general bound.

Example 16. In this example, we compare the performances of two request schemes to motivate the benefit of information optimization. Consider the graph in Figure 6-2 that has n disjoint paths with each path having n edges. We set $\alpha_{ij} = 1$ for all i, j , and let $W_{ij} \sim N(0, 1)$ and independent RVs.

Define the distribution $p_n^s \in \Gamma(N)$ that yields $\hat{W}_{1j} = W_{1j}$ and $\hat{W}_{ij} = 0$ for $i > 1$ (serial information). Define the distribution $p_n^p \in \Gamma(N)$ that yields $\hat{W}_{i1} = W_{i1}$ and $\hat{W}_{ij} = 0$ for $j > 1$ (parallel information).

¹Note that $J(p_{\hat{W}})$ is symmetric with respect to Z and $-Z$, as it should be since Z and $-Z$ have the same distribution.

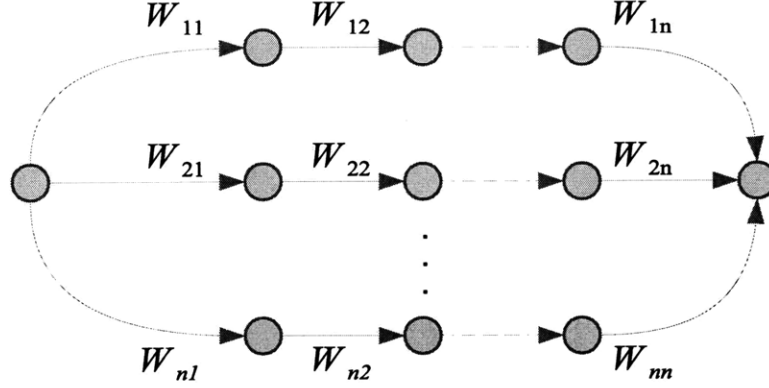


Figure 6-2: Graph for Example 16

The average performances in each case can be written as

$$\begin{aligned}
 J_R(p_n^s) &= \mathbb{E} \left[\min_f \left\{ \left(Rf_1^2 + f_1 \sum_{j=1}^n \hat{W}_{1j} \right) + \sum_{i=2}^n Rf_i^2 \right\} \right] \\
 &= \mathbb{E} \left[\min_f \left\{ f_1 \sqrt{n}Z + R \sum_{i=1}^n f_i^2 \right\} \right] \\
 J_R(p_n^p) &= \mathbb{E} \left[\min_f \left\{ \sum_{i=1}^n Rf_i^2 + \sum_{i=1}^n f_i \hat{W}_{i1} \right\} \right]
 \end{aligned}$$

where $Z \sim N(0, 1)$. It is not clear how to exactly evaluate these expressions, so we instead compare bounds for them instead. We hypothesize that the use of p^s will (at least asymptotically) outperform p^p , so we compare an upper bound for $J_R(p_n^s)$ to a lower bound for $J_R(p_n^p)$.

We derive an upper bound for $J_R(p_n^s)$ using the selection of a suboptimal flow f as well as the inequality $\|f\|_2^2 \leq 1$. Since f will be a function of Z , we choose (the suboptimal) f according to $f_1(Z) = 1$ if $Z < 0$ and $f_1(Z) = 0$ if $Z > 0$. Therefore,

$$J_R(p_n^s) \leq \mathbb{E} \left[\min_f \{ f_1 \sqrt{n}Z + R \} \right] = R + \mathbb{E} \left[\min \{ 0, \sqrt{n}Z \} \right] = R - \sqrt{n} \frac{1}{\sqrt{2\pi}}.$$

We derive a lower bound $J_R(p_n^p)$ by simply splitting the min operator and using the lower bound $\|f\| > 0$.

$$\begin{aligned}
 J_R(p_n^p) &\geq \mathbb{E} \left[\min_f \left\{ \sum_{i=1}^n Rf_i^2 \right\} \right] + \mathbb{E} \left[\min_f \left\{ \sum_{i=1}^n f_i \hat{W}_{i1} \right\} \right] \\
 &\geq 0 + \mathbb{E} \left[\min_i \{ \hat{W}_{i1} \} \right]
 \end{aligned}$$

$$\geq -\sqrt{2 \ln n}$$

where the last inequality is obtained by using Lemma 3 in [8].

Therefore, $J_R(p_n^s)$ (serial information) asymptotically outperforms $J_R(p_n^p)$ (parallel information).

6.4.2 Computational Performance Bounds

We now seek to derive a computational lower bound for performance that is amenable to information optimization. We begin with the follow lemma, which allows us to remove the effect of mean μ from our bounds.

Lemma 5.

$$J_R(p_{\hat{W}}) \geq J(0) - \sum_i \mathbb{E} \left[\min_{f \in \Delta_i} \left\{ \sum_{e \in S_i} f_e \left[\alpha_e R f_e + (\hat{W}_e - \mu_e) \right] \right\} \right]$$

where $J(0) = \min_{p \in P} \{p^T \mu\}$ and Δ_i and f are defined as in Proposition 20.

Proof. Adding and subtracting μ and using the relationship

$$\min_x \{a(x) + b(x)\} \geq \min_x \{a(x)\} + \min_x \{b(x)\}$$

yields

$$\begin{aligned} J_R(p_{\hat{W}}) &\geq \mathbb{E} \left[\min_{f \in \Delta} \left\{ \sum_p f_p \sum_{e \in p} \hat{W}_e \right\} \right] \\ &\quad + \mathbb{E} \left[\min_{f \in \Delta} \left\{ \sum_p f_p \sum_{e \in p} \left[\alpha_e R f'_e + (\hat{W}_e - \mu) \right] \right\} \right] \\ &= \min_{f \in \Delta} \left\{ \sum_p f_p \sum_{e \in p} \mu_e \right\} + \mathbb{E} \left[\min_{f \in \Delta} \left\{ \sum_p f_p \sum_{e \in p} \left[\alpha_e R f'_e + (\hat{W}_e - \mu) \right] \right\} \right] \end{aligned}$$

The first optimization can be interpreted as a congestionless optimization, the minimizing flow for which is the one that has all agents taking the shortest path in the graph. Therefore

$$\geq J(0) + \mathbb{E} \left[\min_f \left\{ \sum_p f_p \sum_{e \in p} \left[\alpha_e R f'_e + (\hat{W}_e - \mu) \right] \right\} \right].$$

Bounding the right-hand term using Proposition 20 proves the claim. \square

We now present the computational lower bound that is the basis for our information optimization algorithm.

Proposition 21.

$$J_R(p_{\hat{W}}) \geq J(0) - \frac{1}{4R} \sum_e \frac{\text{VAR} [\hat{W}_e]}{\alpha_e} - \sum_i \frac{1}{\sqrt{\sum_{e \in S_i} \alpha_e^{-1}}} \sqrt{\sum_{e \in S_i} \frac{\text{VAR} [\hat{W}_e]}{\alpha_e}}.$$

Proof. By Lemma 5, we can compute a lower bound for performance over a single cut and then simply sum over the cuts for the lower bound. Hence, we are interested in computing a lower bound for the expression

$$\min_{f \in \Delta_i} \left\{ \sum_e f_e [\beta_e f_e + (\hat{W}_e - \mu_e)] \right\}$$

where $\beta_e = R\alpha_e$. Assume without loss of generality that $\mu = 0$.

Taking the Lagrangian dual of the optimization, we get

$$= \max_{\gamma, \lambda \geq 0} \left\{ \min_f \left\{ \frac{f^T Q f}{2} + \hat{W}^T f + \gamma 1^T f - 1 - \lambda^T f \right\} \right\},$$

where $Q = \text{diag}\{2\beta_e\}_e$. The optimal f is given by $f = Q^{-1}(\lambda - \gamma 1 - \hat{W})$. Substituting this value for f and simplifying, we get

$$= \frac{1}{2} \max_{\gamma, \lambda \geq 0} \left\{ -(\lambda - \gamma 1 - \hat{W})^T Q^{-1}(\lambda - \gamma 1 - \hat{W}) - 2\gamma \right\},$$

for which the optimizing γ is $\gamma = \frac{1^T Q^{-1}(\lambda - \hat{W}) - 1}{1^T Q^{-1} 1}$. Substituting this value for γ and simplifying, we get

$$= \frac{1}{2} \max_{\lambda \geq 0} \left\{ -(\lambda - \hat{W})^T Q^{-1}(\lambda - \hat{W}) + \frac{1}{1^T Q^{-1} 1} \left[(\lambda - \hat{W})^T Q^{-1} 1 1^T Q^{-1}(\lambda - \hat{W}) - 2(\lambda - \hat{W})^T Q^{-1} 1 + 1 \right] \right\}.$$

We lower bound the expression by removing the positive quadratic term and the positive constant to get

$$= \frac{1}{2} \max_{\lambda \geq 0} \left\{ -(\lambda - \hat{W})^T Q^{-1}(\lambda - \hat{W}) - \frac{2}{1^T Q^{-1} 1} (\lambda - \hat{W})^T Q^{-1} 1 \right\}.$$

We achieve a lower bound for the above expression by choosing a suboptimal λ . The unconstrained minimizing value for λ is given by the solution to the equation

$$-2Q^{-1}(\lambda - \hat{W}) - \frac{2}{1^T Q^{-1} 1} Q^{-1} 1 = 0,$$

for which the solution is $\lambda = -\frac{1}{1^T Q^{-1} 1} + \hat{W}$. However, since λ is constrained to non-negativity, we simply select $\lambda_e = [\hat{W}_e]^+$. Substituting this value for λ and simplifying, we get

$$\geq -\frac{1}{2} \left[-[\hat{W}^T]^- Q^{-1} [\hat{W}]^- + \frac{2}{1^T Q^{-1} 1} [\hat{W}^T]^- Q^{-1} 1 \right].$$

Substituting $Q^{-1} = \frac{1}{2} R^{-1} \cdot \text{diag}\{\alpha_e^{-1}\}_e$, simplifying, and taking expectations, we get

$$= -\frac{1}{4R} \sum_e \frac{\mathbb{E} \left[\left([\hat{W}^T]^- \right)^2 \right]}{\alpha_e} + \frac{2}{\sum_e \alpha_e^{-1}} \sum_e \mathbb{E} [\hat{W}]^- \alpha_e^{-1}.$$

Letting $\gamma_e^2 = \text{VAR} [\hat{W}_e]$, we have the inequalities

$$\begin{aligned} \mathbb{E} \left[\left([\hat{W}^T]^- \right)^2 \right] &\leq \gamma_e^2 \\ \mathbb{E} [\hat{W}]^- &\geq -\frac{1}{2} \sqrt{\gamma_e^2} \end{aligned}$$

where the second inequality comes from the Cauchy-Schwartz inequality [5]. Thus, we have

$$\geq -\frac{1}{4R} \sum_e \frac{\gamma_e^2}{\alpha_e} - \sum_e \frac{\sqrt{\gamma_e^2} \alpha_e^{-1}}{\sum_e \alpha_e^{-1}}.$$

Applying Jensen's Inequality to the second expression gives

$$\geq -\frac{1}{4R} \sum_e \frac{\gamma_e^2}{\alpha_e} - \frac{1}{\sqrt{\sum_e \alpha_e^{-1}}} \sqrt{\sum_e \frac{\gamma_e^2}{\alpha_e}}.$$

Summing over i (the cuts) yields the bound in the claim. \square

6.4.3 Information Optimization

We now present the main result of this section: an information optimization algorithm.

Theorem 11. *A lower bound for $J_R(C)$ is given by the concave maximization*

$$J_R(C) \geq J(0) - \max_{\gamma \geq 0, \|\gamma\|_1 \leq C} \left\{ \frac{1}{4R} \sum_e \frac{\gamma_e}{\alpha_e} + \sum_i \frac{1}{\sqrt{\sum_{e \in S_i} \alpha_e^{-1}}} \sqrt{\sum_{e \in S_i} \frac{\gamma_e}{\alpha_e}} \right\}$$

Proof. Simply substitute $\gamma_e = \text{VAR} [\hat{W}]$ in Proposition 21. \square

6.5 An Analytic Relationship between Capacity and Performance

Finally, we use Proposition 21 to derive an analytic lower bound for $J_R(C)$ in terms of the capacity C , the total flow R , and certain characteristics of the graph when the α_e 's are constant.

Theorem 12. *If $\alpha_e = \alpha$ and $|S_i| \geq M$, then $J_R(C) \geq J(0) - \frac{1}{4\alpha R}C - \sqrt{\frac{N}{M}}\sqrt{C} \sim -O(C)$.*

Proof. Let $\gamma_e^2 = \text{VAR}[\hat{W}_e]$. If $\alpha_e = \alpha$, then by Proposition 21,

$$\begin{aligned} J_R(p_{\hat{W}}) &\geq J(0) - \frac{1}{4\alpha R} \sum_e \gamma_e^2 - \sum_{i=1}^N \sqrt{\frac{\alpha}{|S_i|}} \frac{1}{\sqrt{\alpha}} \sqrt{\sum_{e \in S_i} \gamma_e^2} \\ &\geq J(0) - \frac{1}{4\alpha R} \sum_e \gamma_e^2 - \sqrt{\frac{1}{M}} \sum_{i=1}^N \sqrt{\sum_{e \in S_i} \gamma_e^2}. \end{aligned}$$

Using Jensen's Inequality to add the square root terms:

$$\sum_{i=1}^N \sqrt{X_i} = N \sum_i \frac{\frac{1}{N} \sqrt{X_i}}{\sum_i \frac{1}{N}} \leq N \sqrt{\sum_i \frac{\frac{1}{N} X_i}{\sum_i \frac{1}{N}}} = \sqrt{N} \sqrt{\sum_i X_i}$$

and using the bound $\sum_e \gamma_e^2 \leq C$ yields the bound in the claim. \square

Remark 16. *In shifting the effect of μ into $J(0)$, we get $J_R(0) \geq J(0)$. In general, the two are not equal. Therefore, this bound is necessarily offset from the true performance for any graph where $J_R(0) \neq J(0)$.*

We compare our analytic bound to the actual performance of the optimization for a simple example.

Example 17. *Once again, consider the graph in Figure 6-1. Set $\alpha_1 = \alpha_2 = 1$ and let $W_i \sim N(0, 10)$ and independent. We consider $p_{\hat{W}} \subset \Gamma(C)$ for $C \leq 20$ subject to the additional restriction that \hat{W}_1 and \hat{W}_2 are independent Gaussians as well. We assume a continuous flow of agents, allowing us to simply normalize $R = 1$. One can interpret the use of continuous flow as a normalized approximation to having many agents.*

The plot in Figure 6-3 shows the simulated performance of the agents as the capacity is varied in the range $0 \leq C \leq 20$. The capacity is equally allocated between the two edge weights so that $\text{VAR}[\hat{W}_i] = \sqrt{\frac{C}{2}}$. For each value of C , $J_R(C)$ is approximated by averaging 5,000 simulations.

Interestingly, the performance seems to grow linearly with capacity, as predicted by the bound in Corollary 12.

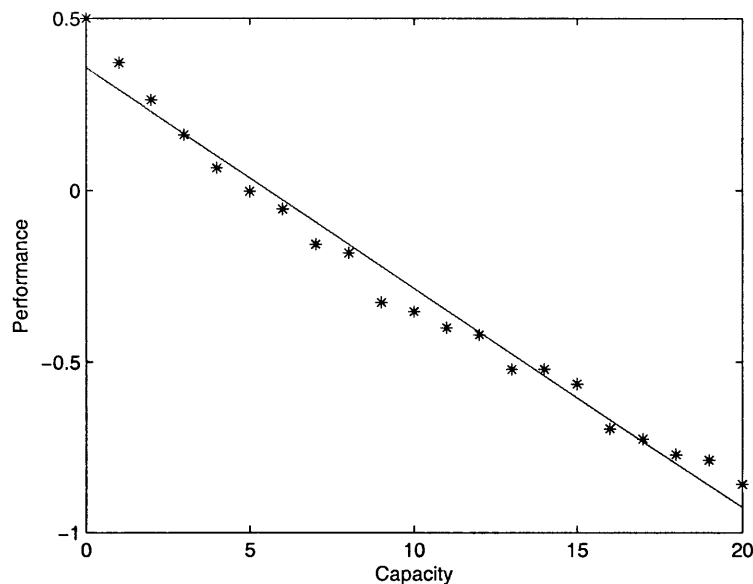


Figure 6-3: Simulated performance of continuous flow with $R = 1$ on the graph in Figure 6-1 as capacity is varied from $0 \leq C \leq 20$. Each point represents a data point from the simulation. The solid line is a linear regression of those points.

6.6 Summary

In this paper, we examined the relationship between information quantity and performance in network flow optimization. It was shown that obtaining analytic bounds for network flow optimization is difficult even in rudimentary cases. Hence, we followed the approach taken Chapter 4 and constructed an outer approximation for the constraint set (based in cuts of the graph) to lower bound the optimization. Both analytic bounds and computational bounds useful for information optimization were developed. The information optimization algorithm in this paper corresponded to a linear optimization that placed more information resources on edges that were more sensitive to congestion. Overall, the fastest rate of improvement with information is linear in the amount of information available to the agents. An example illustrated that linear growth with capacity may be a tight bound on growth.

Chapter 7

Conclusions and Future Work

In this thesis, we examined the value of information in an instance of an uncertain decision framework: shortest path optimization on a graph with random edge weights. In the most basic version of the problem, we considered an agent receiving side information about the edge weights in advance of its travel, using that information to estimate the edge weights, and then traveling along the path with the smallest estimated length. We sought to understand how side information impacted the agent’s average performance as well as to understand how the agent should optimize its side information.

We defined a notion of information that was compatible with this formulation: the variance of the edge-weight estimates. By allowing the agent to optimize the information it receives, we needed to develop algorithms for information optimization. A practical set of such algorithms was developed by allowing the agent to optimize bounds for average performance. The bounds were based in a new graph reduction for shortest path optimization that simplified the description of the graph to the minimal-sphere and subspace containing it. In certain cases, the optimization can be practically bounded by a convex optimization.

We also provided a generalization to a sequential-information framework whereby the agent can receive and apply information as it traverses the graph. We considered two cases of sequential information, controlled and uncontrolled. In the controlled case, we allowed the agent to optimize information it would receive at each stage at that stage. In the uncontrolled case, we forced the agent to optimize the information it would receive at each stage before traversing the graph altogether. We defined a new information constraint set for this case (though still based in the variance of the edge-weight estimate) and developed a new abstraction for sequential decision making that greatly simplified analysis. The lower bounds we developed were a generalization of those in the non-sequential case, and the bounds indicated a loss of performance in the uncontrolled-information case while showing a potential gain in performance in the controlled-information case.

Finally, we applied our framework the study the value of information in network flow optimization. It was shown that obtaining analytic performance bounds for network flow optimization is difficult even in simple cases, and, hence, we followed the approach taken in shortest path optimization by constructing an outer approximation for the constraint set

(based in cuts of the graph, not spheres). Both analytic bounds and computational bounds useful for information optimization were developed.

Future work includes further examining the relationship between spheric outer approximation and the edge-cut outer approximations used in Chapters 4 and 6 to actual rates of improvement. The bounds seem to capture the best possible rates of growth, and relating them to optimizations such as that in Corollary 7 may reveal important general structures that can be used to bound more general partial-information optimizations. Future work also includes better understanding the impact of controlled information in the sequential information setting by developing bounds that leverage coarser graph topology information and developing broader conditions under which controlled sequential information can outperform non-sequential information.

Appendix A

Additional Proofs and Results

Chapter 4

Proof of Theorem 2. We prove the result by providing such an algorithm. First, define

$$\begin{aligned}
 f(\{n^1, \dots, n^i\}, j) &= \underset{(n, \Delta)}{\operatorname{argmin}} \|n\| \text{ subject to} \\
 n_j &= 1, \quad (n^k)^T n = 0 \text{ for all } k \leq i \\
 \overline{J}(v) &= \max_{e \mid \operatorname{hd}(e)=v} \{\overline{J}(\operatorname{tl}(e)) + n_e\}, \quad \underline{J}(v) = \min_{e \mid \operatorname{hd}(e)=v} \{\underline{J}(\operatorname{tl}(e)) + n_e\} \\
 \overline{J}(s) &= \underline{J}(s), \quad \overline{J}(t) = \underline{J}(t) = 0
 \end{aligned}$$

with the definition $f(\{n^1, \dots, n^{i-1}\}, j) = (0, \infty)$ if the optimization is infeasible.

The algorithm is as follows:

```

1:  $H_{\mathcal{P}} \leftarrow I; i \leftarrow 1; j \leftarrow 1$ 
2: while  $j \leq |E|$  do
3:    $(n^*, \Delta^*) \leftarrow f(\{n^1, \dots, n^{i-1}\}, j)$ 
4:   if  $\Delta^* = 0$  then
5:      $H_{\mathcal{P}} \leftarrow H_{\mathcal{P}} - \frac{n^*(n^*)^T}{\|n^*\|^2}$ 
6:      $n^i \leftarrow n^*$ 
7:      $i \leftarrow i + 1$ 
8:   else
9:      $j \leftarrow j + 1$ 
10:  end if
11: end while

```

It is clear that the above algorithm terminates in polynomial number of executions, so we only need to prove that it generates a valid projection matrix for $S_{\mathcal{P}}$. For any orthogonal basis $\{n^1, n^2, \dots, n^k\}$ of $S_{\mathcal{P}}^\perp$, a projection matrix for $S_{\mathcal{P}}$ can be written as $I - \sum_{i=1}^m \frac{n^i(n^i)^T}{\|n^i\|^2}$. Therefore, we need to show that the set $\{n^i\}$ generated from the algorithm is such a basis.

A vector $n \in \mathbb{R}^{|E|}$ being normal to $S_{\mathcal{P}}$ is equivalent to $n^T (\mathcal{P} - \overline{p}) = \{0\} \Leftrightarrow n^T (p - \overline{p}) = 0$

for all $p \in P$. Therefore, a sufficient and necessary condition for n to be normal is that all paths from s to t have the same length when the edge weights are n .

The longest path in the graph when the edge weights are given by n is the unique function $\bar{J} : V \rightarrow \mathbb{R}$ satisfying $\bar{J}(v) = \max_{e \mid \text{hd}(e)=v} \{\bar{J}(\text{tl}(e)) + n_e\}$. The shortest path is the unique function \underline{J} satisfying the same equality with $\max \rightarrow \min$. An equivalent condition to all paths having the same length is $\bar{J}(s) = \underline{J}(s)$. Therefore, if f is feasible, any optimal solution vector n^* provided by f must be normal to $S_{\mathcal{P}}$, and, further, it must be non-zero. Therefore, any set of vectors $\{n^1, n^2, \dots, n^k\}$ must lie in $S_{\mathcal{P}}^\perp$.

By the orthogonality constraint in f , any set of vectors $\{n^1, n^2, \dots, n^k\}$ yielded by the algorithm must be non-zero and orthogonal, hence the set is a subset of an orthogonal basis for $S_{\mathcal{P}}^\perp$.

Now, suppose the set $\{n^1, n^2, \dots, n^k\}$ is a strict subset of a basis. Then there is a non-zero vector n that is orthogonal to each n^j and lies in $S_{\mathcal{P}}^\perp$. Let i be the smallest non-zero component of n , and assume without loss of generality that $n_i = 1$ (n can be normalized to produce this). Finally, let $l = \arg\max_j \{n^j | n_i^j = 1\}$.

By the definition of l , $f(\{n^1, \dots, n^l\}, i)$ is infeasible since, otherwise, $n_i^{l+1} = 1$ yielding $l = l + 1$. By the existence of n , though, we know that $f(\{n^1, \dots, n^l\}, i)$ is feasible since n satisfies all of the constraints of the optimization. The contradiction implies that $\{n^1, n^2, \dots, n^k\}$ must be an orthogonal basis for $S_{\mathcal{P}}^\perp$. Hence, $H_{\mathcal{P}}$ is a projection matrix for $S_{\mathcal{P}}$. \square

Proof of Corollary 7. We present a detailed sketch of the proof. First, if we remove the capacity constraint from (4.3) and instead fix the variances $\text{VAR}[\hat{W}_e] = \lambda_e^2$, get

$$\underline{J}(\{\lambda_e\}) = \min_{P_W} \left\{ \mathbb{E} \left[\min_p \{p^T \hat{W}\} \right] \right\} \text{ subject to } \mathbb{E}[\hat{W}] = \mu, \text{ VAR}[\hat{W}] = \lambda_e^2$$

where, for ease, we denote the lower bound for performance in this case as $J(\lambda)$. This optimization is of the form Equation (3.7) in [4]. By Theorem 3.1 in [4], it is equivalent to Equation (3.8) in [4]. Substituting our constraints yields a quadratic objective and second-degree polynomial inequalities, which we can re-express as operations and inequalities on semi-definite matrices:

$$\underline{J}(\{\lambda_e\}) = \max_{\{G_e\}, d \in \mathbb{R}^{|E|}} \left\{ J(d) + \sum_e G_e \cdot \begin{bmatrix} \lambda_e^2 + \mu_e^2 & \mu_e \\ \mu_e & 1 \end{bmatrix} \right\} \text{ subject to } G_e \leq \left\{ 0, \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & -d_e \end{bmatrix} \right\}$$

where by our definition of $J(W)$, $J(d)$ is simply $J(W)$ with $W = d$, a constant.

Let $\Gamma'(C) = \{\lambda \in \mathbb{R}^{|E|} \mid 0 \leq \lambda_e^2 \leq \sigma_e^2, \|\lambda\|_2^2 \leq C\}$. $\Gamma'(C)$ is clearly a convex set, and any

$p_{\hat{W}} \in \Gamma(C)$ will yield a set of variances $\text{VAR}[\hat{W}_e] = \lambda_e^2$ in $\Gamma'(C)$. We have

$$\underline{J}(C) = \min_{\lambda \in \Gamma'(C)} \underline{J}(\lambda).$$

Taking the dual of the inner optimization for $\underline{J}(\gamma)$ yields a new inner optimization

$$\begin{aligned} \min_{H_e} \max_d \left\{ J(d) - \sum_e H_e \cdot \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & -d_e \end{bmatrix} \right\} \text{ subject to} \\ H_e \leq 0, \quad H_e \geq - \begin{bmatrix} \lambda_e^2 + \mu_e^2 & \mu_e \\ \mu_e & 1 \end{bmatrix}. \end{aligned}$$

Let $H_e = \begin{bmatrix} a_e & b_e \\ b_e & c_e \end{bmatrix}$. Letting $c = [c_1 \dots c_{|E|}]^T$ and $d = [d_1 \dots d_{|E|}]^T$, the objective in the minimax is

$$\begin{aligned} &= J(d) + c^T d - \sum_e \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{bmatrix} \cdot H_e \\ &= \min_{p \in \mathcal{P}} \{p^T d\} + c^T d - \sum_e \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{bmatrix} \cdot H_e \\ &= \min_{p \in \mathcal{P}} \{(p + c)^T d\} + \sum_e \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix} \cdot H_e \end{aligned}$$

Of interest to us is the minimaximin expression:

$$\min_c \max_d \min_{p \in \mathcal{P}} \{(p + c)^T d\}$$

If $-c \in \mathcal{P}$, then this expression must always be nonpositive since $0 \in \mathcal{P} - c$. If $-c \notin \mathcal{P}$, then one can show that the expression will always be ∞ . Therefore, we require $-c \in \mathcal{P}$. In this case, $d = 0$ is the optimal strategy for d since that maximizes the expression to 0.

The constraint $c \in \mathcal{P}$ is represented by

$$\sum_e \left(\begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix} \cdot H_e \right) v_e \in \mathcal{P}.$$

The remainder of the claim easily follows. □

Chapter 5

Proposition 22.

$$\begin{aligned} \mathbb{E} [J_n(Y_n, Y_{n-1}, S_{n-1}) | Y_{n-1}] &\geq -\sqrt{(r_o^*)^2 - \gamma_{n-1}^2} \left(\mathbb{E} [\hat{W}_n^T H_{n-1} \hat{W}_n | Y_{n-1}] \right. \\ &\quad \left. + 2 \mathbb{E} [\hat{W}_n^T H_{n-1} \hat{W}_{n-1} | Y_{n-1}] + \hat{W}_{n-1}^T H_{n-1} \hat{W}_{n-1} \right)^{\frac{1}{2}} \\ &\quad + c_{n-1} \hat{W}_{n-1}. \end{aligned}$$

Proof. We start with the lower bound from Lemma 4, and we assume without loss of generality that $\mu = 0$ (so that $J(0) = 0$ as well). We first “center” $\overline{\mathcal{P}} \cap S$ by removing c_o^* . First, we have

$$\begin{aligned} \underline{J}_1(Y_1) &= J(0) + \min_{S_1 \in \mathcal{S}_1} \{ \mathbb{E} [\underline{J}_2(Y_2, Y_1, S_1) | Y_1] \} \\ \underline{J}_i(Y_i, Y_{i-1}, S_{i-1}) &= \min_{S_i \in \mathcal{S}_i(S_{i-1})} \{ \mathbb{E} [\underline{J}_{i+1}(Y_{i+1}, Y_i, S_i) | Y_i] \} \\ \underline{J}_n(Y_n, Y_{n-1}, S_{n-1}) &= \min_{p \in (\overline{\mathcal{P}} - c_o^*) \cap S_{n-1}} \{ p^T \hat{W}_{\overline{n}} \} + (c_o^*)^T \hat{W}_{\overline{n}}. \end{aligned}$$

It is straightforward that $\underline{J}_1 = \underline{J}'_1$ where

$$\begin{aligned} \underline{J}'_1(Y_1) &= \min_{S_1 \in \mathcal{S}_1} \{ \mathbb{E} [\underline{J}_2(Y_2, Y_1, S_1) | Y_1] \} + \mathbb{E} [(c_o^*)^T \hat{W}_{\overline{n}}] \\ \underline{J}'_i(Y_i, Y_{i-1}, S_{i-1}) &= \min_{S_i \in \mathcal{S}_i(S_{i-1})} \{ \mathbb{E} [\underline{J}_{i+1}(Y_{i+1}, Y_i, S_i) | Y_i] \} \\ \underline{J}'_n(Y_n, Y_{n-1}, S_{n-1}) &= \min_{p \in (\overline{\mathcal{P}} - c_o^*) \cap S_{n-1}} \{ p^T \hat{W}_{\overline{n}} \}. \end{aligned}$$

where $\mathbb{E} [(c_o^*)^T \hat{W}_{\overline{n}}] = 0$. Therefore, assume without loss of generality that $c_o^* = 0$. Note that this implies that $\overline{\mathcal{P}} = B(0, r_o^*)$.

By definition, $\overline{\mathcal{P}} \cap S_{n-1} = B(c_{n-1}, r') \cap S_{n-1}$ for some $r' \leq r$. For ease, call this intersection B . Since $B \subset S_{n-1}$, $H_{n-1}B = B$. Therefore,

$$\begin{aligned} \underline{J}_n(Y_n, Y_{n-1}, S_{n-1}) &= \min_{p \in B} \{ p^T \hat{W}_{\overline{n}} \} \\ &= \min_{p \in (B - c_{n-1})} \{ p^T \hat{W}_{\overline{n}} \} + c_{n-1}^T \hat{W}_{\overline{n}} \\ &= \min_{p \in (B - c_{n-1})} \{ (H_{n-1}p)^T \hat{W}_{\overline{n}} \} + c_{n-1}^T \hat{W}_{\overline{n}} \\ &= -r' \|H_{n-1} \hat{W}_{\overline{n}}\| + c_{n-1}^T \hat{W}_{\overline{n}}. \end{aligned}$$

By the relationship

$$\mathbb{E} \left[\hat{W}_{\bar{j}} | Y_{\bar{j}-1} \right] = \mathbb{E} \left[\mathbb{E} \left[W | Y_{\bar{j}} \right] | Y_{\bar{j}-1} \right] = \mathbb{E} \left[W | Y_{\bar{j}-1} \right] = \hat{W}_{\bar{j}-1}$$

and Jensen's Inequality, we have

$$\mathbb{E} \left[\underline{J}_n(Y_n, Y_{n-1}, S_{n-1}) | Y_{n-1} \right] = -r' \sqrt{\mathbb{E} \left[\|H_{n-1} \hat{W}_{\bar{n}}\|^2 | Y_{n-1} \right]} + c_{n-1}^T \hat{W}_{n-1}.$$

Using the relationship $\hat{W}_{\bar{n}} = \hat{W}_n + \hat{W}_{n-1}$ and $H_{n-1}H_{n-1} = H_{n-1}$, we have

$$\|H_{n-1} \hat{W}_{\bar{n}}\|^2 = \hat{W}_n^T H_{n-1} \hat{W}_n + 2\hat{W}_n^T H_{n-1} \hat{W}_{n-1} + \hat{W}_{n-1}^T H_{n-1} \hat{W}_{n-1}.$$

Finally, it is straightforward to show that $(r')^2 = (r_o^*)^2 - \|c_{n-1}\|^2 = (r_o^*)^2 - \gamma_{n-1}^2$. The bound in the claim follows from the fact that $J_n \geq \underline{J}_n$. \square

Appendix B

The Value of Information in Energy Production

We consider an application of the general framework presented in Chapter 2 to energy production whereby the energy supplier seeks to minimize the cost of energy production subject to having some information about (a random) demand. Demand information is important to the supplier in our setting because the supplier is penalized if it undersupplies demand.

Let $D(t)$ be the random demand for power at time t with $E[D(t)] = \mu(t)$ and $\text{VAR}[D(t)] = \sigma^2(T)$, and let $p(t)$ be the power produced at time t . The supplier's cost $c(p, D)$ of production is a function of both p and D , which, in our simplified setting, is defined as

$$c(p, D) = \sum_t \alpha p(t) + \sum_t \beta [D(t) - p(t)]^+$$

where α is the unit cost of power, and β is the unit cost (in the form of a penalty, such as a fine) of undersupplying demand.

We have some side information Y about D , and we assume that the side information has a certain amount of capacity $C(t)$ at each time t . We define our information set $\Gamma(C)$ similarly as before: $\Gamma(C) = \{p_{DY} \mid \text{VAR}[\hat{D}(t)] \leq C(t)\}$.

We seek to determine the value of demand information to the supplier. Hence, we seek to bound

$$\begin{aligned} J(C) &= \min_{p_{DY} \in \Gamma(C)} \left\{ E \left[\min_p \{E[c(p, D)|Y]\} \right] \right\} \\ &= \min_{p_{DY} \in \Gamma(C)} \left\{ E \left[\min_p \left\{ \sum_t \alpha p(t) + \sum_t \beta E[[D(t) - p(t)]^+ | Y] \right\} \right] \right\}. \end{aligned}$$

We upper bound this quantity by choosing a suboptimal value for $p(t)$ and applying

Jensen's Inequality. We first reparameterize the inner optimization using $\hat{p} = p - \hat{D}$.

$$\begin{aligned} & \mathbb{E} \left[\min_{\hat{p}} \left\{ \sum_t \alpha p(t) + \sum_t \beta \mathbb{E} [[D(t) - p(t)]^+ | Y] \right\} \right] \\ &= \mathbb{E} \left[\min_{\hat{p}} \left\{ \sum_t \alpha (\hat{p}(t) + \hat{D}(t)) + \sum_t \beta \mathbb{E} [[D(t) - \hat{D}(t) - \hat{p}(t)]^+ | Y] \right\} \right]. \end{aligned}$$

Applying Jensen's Inequality yields

$$\leq \min_{\hat{p}} \left\{ \mathbb{E} \left[\sum_t \alpha (\hat{p}(t) + \hat{D}(t)) \right] + \sum_t \beta \mathbb{E} [[D(t) - \hat{D}(t) - \hat{p}(t)]^+] \right\}.$$

Finally, choosing a suboptimal solution $\hat{p} = 0$ yields an upper bound strictly in terms of capacity:

$$J(C) \leq \alpha \sum_t \mu(t) + \frac{\beta}{2} \sum_t \sqrt{\sigma^2(t) - C(t)}$$

where the inequality $\mathbb{E} [D(t) - \hat{D}(t)]^+ \leq \frac{1}{2} \sqrt{\text{VAR} [D(t) - \hat{D}(t)]}$ follows from [4]

Appendix C

Upper Bounds via Splitting and Pruning

C.1 Computing an Upper Bound by Splitting the Graph

We now consider an alternative information optimization that optimizes an upper bound for $J(p_{\hat{W}})$ that better leverages the structure of the graph to reduce complexity. The upper bound is based in Jensen's Inequality, and it applies a dynamic-programming-like approach to computing the average performance backward through the graph, but it requires that the edge weights obey the Gaussian assumptions discussed before.

Theorem 13. *An upper bound $J(C)$ under independent Gaussian edge weights is*

$$J(C) \leq \min_{\{\gamma_e^2\} \in \Gamma(C)} \{\bar{J}(s, \{\gamma_e^2\})\} \quad (\text{C.1})$$

where

$$\bar{J}(v, \{\gamma_e^2\}) = \int j \frac{\partial}{\partial x} \left(\prod_{e \mid \text{hd}(e)=v} \left[1 - \Phi \left(\frac{x - \bar{J}(\text{tl}(e), \{\gamma_e^2\}) - \mu_e}{\gamma_e} \right) \right] \right) \Big|_{x=j} dj$$

and $\bar{J}(t, \{\gamma_e^2\}) = 0$.

Proof. Let $J(v, \hat{W})$ be the length of the shortest path from vertex v to t under edge weights \hat{W} :

$$J(v, \hat{W}) = \min_{e \mid \text{hd}(e)=v} \left\{ J(\text{tl}(e), \hat{W}) + \hat{W}_e \right\},$$

$$J(t, \hat{W}) = 0.$$

By Jensen's Inequality and acyclicity, we have

$$\mathbb{E} [J(v, \hat{W})] \leq \mathbb{E} \left[\min_{e | \text{hd}(e)=v} \left\{ \mathbb{E} [J(\text{tl}(e), \hat{W})] + \hat{W}_e \right\} \right],$$

so that the set of equations

$$\bar{J}(v, \{\gamma_e^2\}) = \mathbb{E} \left[\min_{e | \text{hd}(e)=v} \left\{ \bar{J}(\text{tl}(e), \{\gamma_e^2\}) + \hat{W}_e \right\} \right]$$

yields an upper bound for $\mathbb{E} [J(v, \hat{W})]$.

Now let,

$$\bar{J}(v, \{\gamma_e^2\}, \hat{W}) = \min_{e | \text{hd}(e)=v} \left\{ \bar{J}(\text{tl}(e), \{\gamma_e^2\}) + \hat{W}_e \right\}$$

so that $\bar{J}(v, \{\gamma_e^2\}) = \mathbb{E} [\bar{J}(v, \{\gamma_e^2\}, \hat{W})]$.

Using

$$P(\hat{W}_e > x) = 1 - \Phi \left(\frac{x - \mu_e}{\gamma_e} \right)$$

and $P(\min_i \{X_i\} > x) = \prod_i P(X_i > x)$ for independent RVs $\{X_i\}$ allows us to compute the CDF for the $\bar{J}(v, \{\gamma_e^2\}, \hat{W})$. Taking a derivative and integrating yields its expected value:

$$\bar{J}(v, \{\gamma_e^2\}) = \int j \frac{\partial}{\partial x} \left(\prod_{e | \text{hd}(e)=v} \left[1 - \Phi \left(\frac{x - \bar{J}(\text{tl}(e), \{\gamma_e^2\}) - \mu_e}{\gamma_e} \right) \right] \right) \Big|_{x=j} dj$$

□

C.2 Upper Bound on Lost Performance from Pruning the Graph

The information optimizations presented in this thesis are quite manageable, but we can further improve their performances by leveraging a simple fact about real-world graphs: paths of long average length are almost never the shortest path, so information about them can be neglected. We now provide a bound on the performance lost from pruning such paths. In the framework of optimizing over the path polytope, this corresponds to removing certain extreme points from \mathcal{P} .

Once again, define the function $\Theta_X(c) = \mathbb{E} [\min \{X, c\}]$. If $X \sim N(0, 1)$, then $\Theta_X(0) = \frac{-1}{\sqrt{2\pi}}$. We use Θ to write a very simple pruning algorithm. Let p^* be the path in the graph with shortest average length. We want to prune the paths p that are often not the shortest path when compared to p^* . Formally, we prune every path p satisfying

$$\mathbb{E} [\min \{(p^*)^T W, p^T W\}] - (p^*)^T \mu = \Theta_{(p-p^*)^T W}(0) \approx 0.$$

If $E[p^T W]$ is large, we expect that we should be able to prune p , but without additional assumptions on the relationship between the mean $E[p^T W]$ and variance $\text{VAR}[p^T W]$, it may be the case that p is not pruned, even for very large mean.

C.2.1 Efficient Pruning in the Gaussian Case

We seek to establish easily-verifiable conditions under which can efficiently prune paths as well as compute (or estimate) the loss in performance incurred from the pruning. Although the results in this section are specific to the Gaussian case using the information set $\Gamma_G(C)$, they are very easily generalizable to the case of general distributions.

- $\lim_{c \rightarrow \infty} \Theta_X(c) = E[X]$,
- $\lim_{c \rightarrow -\infty} \Theta_X(c) = c$,
- and $\frac{\partial \Theta_X}{\partial c}(c) = 1 - F(c)$ with F being the CDF for X .

We can also show that, under certain conditions, as the variance of X increases, the value of $\Theta_X(c)$ decreases.

Proposition 23. *Let $Z \sim N(0, 1)$, and let $X_1 \stackrel{d}{=} \sigma_1 Z + \mu$ and $X_2 \stackrel{d}{=} \sigma_2 Z + \mu$. If $c \leq \mu$ and $\sigma_1 \leq \sigma_2$, then $\Theta_{X_2}(c) \leq \Theta_{X_1}(c)$.*

Proof. First, notice that

$$\Theta_{X_i}(c) = \mu + \sigma_i \Theta\left(\frac{c - \mu}{\sigma_i}\right).$$

Now, taking a derivative with respect to σ , we get

$$\frac{\partial}{\partial \sigma} \left\{ \mu + \sigma \Theta\left(\frac{c - \mu}{\sigma}\right) \right\} = \left[\Theta\left(\frac{c - \mu}{\sigma}\right) - \frac{c - \mu}{\sigma} \right] + \left[F\left(\frac{c - \mu}{\sigma}\right) - 1 \right] + F\left(\frac{c - \mu}{\sigma}\right) \frac{c - \mu}{\sigma},$$

which is non-positive if $c \leq \mu$. Since

$$\int_{\sigma_1}^{\sigma_2} \frac{\partial}{\partial \sigma} \left(\mu + \sigma \Theta\left(\frac{c - \mu}{\sigma}\right) \right) = \Theta_{X_2}(c) - \Theta_{X_1}(c),$$

the difference is non-positive if the integrand is non-positive. \square

We can now we establish a useful relationship between the mean and variance of a path that will guarantee that as the mean of the path increases, the path's performance relative to \bar{p} becomes inconsequential.

Proposition 24. *Let $Z \sim N(0, 1)$, and let $\{X_\mu\}_\mu$ be a family of RVs indexed by μ with $X_\mu \stackrel{d}{=} \sigma(\mu)Z + \mu$. If $\sigma^2(\mu) \leq \mu$, then*

$$\lim_{\mu \rightarrow \infty} \Theta_{X_\mu}(c) = c.$$

Proof. Assume $\mu \geq c$. By the assumption that $\sigma_\mu \leq \sqrt{\mu}$ and Proposition 23, we have

$$c \geq \Theta_{X_\mu}(c) = \Theta_{\sigma(\mu)Z+\mu}(c) \geq \Theta_{\sqrt{\mu}Z+\mu}(c),$$

which implies

$$|c - \Theta_{\sqrt{\mu}Z+\mu}(c)| \geq |c - \Theta_{X_\mu}(c)|.$$

Therefore, we only need to prove convergence of the lefthand term.

Now,

$$\lim_{\mu \rightarrow \infty} c - \Theta_{\sqrt{\mu}Z+\mu}(c) = \lim_{\mu \rightarrow \infty} \sqrt{\mu} \left(\frac{c - \mu}{\sqrt{\mu}} - \Theta \left(\frac{c - \mu}{\sqrt{\mu}} \right) \right).$$

As $\mu \rightarrow \infty$, the inner term goes to 0, and so we apply L'Hopital's rule to get

$$= \lim_{\mu \rightarrow \infty} \Phi \left(\frac{c - \mu}{\sqrt{\mu}} \right) (1 - 2\mu).$$

Applying L'Hopital's rule once again yields

$$= \lim_{\mu \rightarrow \infty} \phi \left(\frac{c - \mu}{\sqrt{\mu}} \right) \left[\frac{1}{2} \mu^{\frac{3}{2}} - c\sqrt{\mu} \right].$$

Since $\phi(x)x^{\frac{3}{2}} \rightarrow 0$, the limit converges to 0. \square

Verifying that $\mathbb{E}[p^T W]$ and $\text{VAR}[p^T W]$ satisfy the conditions of Lemma 24 for every path p can be cumbersome in general, but there is an easily-verifiable sufficient condition for it.

Proposition 25. *If each edge weight W_e in G satisfies $\sigma_e \leq \sqrt{\mu_e}$, then for each path P in G , $\sigma_P \leq \sqrt{\mu_P}$.*

Proof. For a path P , we have

$$\sigma_P = \sqrt{\sum_{e \in P} \sigma_e^2} \leq \sqrt{\sum_{e \in P} \mu_e} = \sqrt{\mu_P}.$$

\square

Remark 17. *Neither Proposition 23, Lemma 24, nor Proposition 25 really rely on $Z \sim N(0, 1)$, but rather that $\mathbb{E}[Z] = 0$ and $\text{VAR}[Z] = 1$. The Gaussian restriction in this section is because we are assuming $p_{\hat{W}} \in \Gamma_G(C)$.*

We can use Proposition 25 to derive a simple, efficient pruning algorithm:

- 1: Choose some threshold length L .
- 2: For each vertex $v \in V$, compute the average shortest path length $\underline{J}(v)$ from v to t .
- 3: Let $\tilde{V} = \{v \in V \text{ such that } \underline{J}(v) \leq L\}$.
- 4: Let $\tilde{E} = \{e \in E \text{ such that } \text{hd}(e), \text{tl}(e) \in \tilde{V}\}$.
- 5: Construct a new graph $\tilde{G} = (\tilde{V}, \tilde{E})$.

C.2.2 Performance Loss from Path Pruning

The following theorem provides an upper bound for the performance lost from pruning.

Theorem 14. *Let $\mathcal{P} = \bar{\mathcal{P}} \cup \underline{\mathcal{P}}$ where $\bar{\mathcal{P}} = \{p \mid \mathbb{E}[p^T W] \geq L\}$ and $\underline{\mathcal{P}} = \{p \mid \mathbb{E}[p^T W] < L\}$. Then for any $p_{\hat{W}} \in \Gamma_G(C)$,*

$$0 \leq \mathbb{E} \left[\min_{p \in \underline{\mathcal{P}}} \{p^T \hat{W}\} \right] - J(p_{\hat{W}}) \leq \sum_{k > L} \left[(\# \text{ paths of length } k) \left((\bar{L} - k) - \sqrt{k} \Theta \left(\frac{\bar{L} - k}{\sqrt{k}} \right) \right) \right],$$

where $\bar{L} = \bar{p}^T \mathbb{E}[W]$ and $\lim_{k \rightarrow \infty} \left((\bar{L} - k) - \sqrt{k} \Theta \left(\frac{\bar{L} - k}{\sqrt{k}} \right) \right) = 0$.

Proof. The left inequality is clear, so we proceed to prove the right inequality. First,

$$\min_{p \in \mathcal{P}} \{p^T \hat{W}\} = \min_{p \in \bar{\mathcal{P}} \cup \underline{\mathcal{P}}} \{p^T \hat{W}\} = \bar{p}^T \hat{W} + \min_{p \in \bar{\mathcal{P}} \cup \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}\}.$$

Since $0 \in \{(p - \bar{p})^T W\}$ (take $p = \bar{p}$),

$$\begin{aligned} &= \bar{p}^T \hat{W} + \min \left\{ \min_{p \in \bar{\mathcal{P}} \cup \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \right\} \\ &= \bar{p}^T \hat{W} + \min \left\{ \min_{p \in \bar{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\}, \min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \right\} \\ &\geq \bar{p}^T \hat{W} + \min \left\{ \min_{p \in \bar{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \right\} + \min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \\ &= \bar{p}^T \hat{W} + \min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} + \min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}\}, \end{aligned}$$

where the inequality comes from the fact that $\min\{a, b\} \geq a + b$ for $a, b \leq 0$, and the last equality comes from the fact that $\bar{p} \in \underline{\mathcal{P}}$.

Therefore,

$$\begin{aligned} &\min_{\underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}\} - \min_{p \in \mathcal{P}} \{p^T \hat{W}\} \leq -\bar{p}^T \hat{W} - \min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \\ \Rightarrow &\min_{\underline{\mathcal{P}}} \{p^T \hat{W}\} - \min_{p \in \mathcal{P}} \{p^T \hat{W}\} \leq -\min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\}. \end{aligned}$$

Take expectations of both sides. We can upper bound the right side as

$$-\mathbb{E} \left[\min_{p \in \underline{\mathcal{P}}} \{(p - \bar{p})^T \hat{W}, 0\} \right] \leq -\sum_p \mathbb{E} \left[\min \{(p - \bar{p})^T \hat{W}\}, 0 \right] = -\sum_p \Theta_{(p - \bar{p})^T \hat{W}}(0).$$

Let $\sigma_p = \sqrt{\text{VAR}[(p - \bar{p})^T \hat{W}]}$ and let $\mu_p = (p - \bar{p})^T \mathbb{E}[W]$. By the assumption that $p_{\hat{W}} \in \Gamma_G(C)$, we have $(p - \bar{p})^T \hat{W} \stackrel{d}{=} \sigma_p Z_p + \mu_p$ for $Z_p \sim N(0, 1)$.

By Proposition 23, we can upper bound $-\Theta_{(p-\bar{p})^T \hat{W}}(0)$ by maximizing the variance of $\sigma_p = \sqrt{\mathbb{E}[p^T W]}$. Therefore, an upper bound for the right hand side is

$$-\sum_p \left[\mu_p + \sqrt{\mathbb{E}[p^T W]} \Theta \left(\frac{-\mu_p}{\sqrt{\mathbb{E}[p^T W]}} \right) \right].$$

We now simply group together the paths by their lengths, yielding

$$\sum_{k>L} \left[(\# \text{ paths of length } k) \left((\bar{L} - k) - \sqrt{k} \Theta \left(\frac{\bar{L} - k}{\sqrt{k}} \right) \right) \right].$$

Finally, the last claim concerning the limit to zero is simply a restatement of Proposition 24 □

Remark 18. *While the bound in Theorem 14 may not be practical for general computation, it may be applicable to classes of graphs that possess certain regular structures. Noteworthy in the bound is that although the number of paths of a particular long length may be large, Proposition 24 tells us that as long as the number of such paths does not increase too quickly, the performance loss from ignoring them is bounded.*

Bibliography

- [1] Terje Aven. Upper (lower) bounds on the mean of the maximum (minimum) of a number of random variables. *Journal of applied probability*, 22(3):723–728, 1985.
- [2] M. Avriel and A. C. Williams. The value of information and stochastic programming. *Operations Research*, 18(5):947–954, 1970.
- [3] Aharon Ben-Tal and Eithan Hochman. Stochastic programs with incomplete information. *Operations Research*, 24(2):336–347, 1976.
- [4] Dimitris Bertsimas, Karthik Natarajan, and Chung-Piaw Teo. Probabilistic combinatorial optimization: Moments, semidefinite programming, and asymptotic bounds. *SIAM J. on Optimization*, 15(1):185–209, 2005.
- [5] Dimitris Bertsimas, Karthik Natarajan, and Chung-Piaw Teo. Tight bounds on expected order statistics. *Probab. Eng. Inf. Sci.*, 20(4):667–686, 2006.
- [6] Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [7] Fan Chung and Linyuan Lu. The diameter of sparse random graphs. *Adv. Appl. Math.*, 26(4):257–279, 2001.
- [8] Randy Cogill. Randomized load balancing with non-uniform task lengths. In *Forty-Fifth Annual Allerton Conference*, September 2007.
- [9] Luc P. Devroye. Inequalities for the completion times of stochastic pert networks. *Mathematics of Operations Research*, 4(4):441–447, 1979.
- [10] V. Fenchau. Bounds for the expected value of information. *OR Spectrum*, 2(2):65–73, June 1980.
- [11] M. Grant and S. Boyd. Cvx: Matlab software for disciplined convex programming (web page and software). <http://stanford.edu/~boyd/cvx>, June 2009.
- [12] J. H. Hagstrom. Computational complexity of pert problems. *Networks*, 18:47–56, 1988.

- [13] Refael Hassin and Eitan Zemel. On shortest paths in graphs with random weights. *Mathematics of Operations Research*, 10(4):557–564, 1985.
- [14] C. C. Huang, I. Vertinsky, and W. T. Ziemba. Sharp bounds on the value of perfect information. *Operations Research*, 25(1):128–139, 1977.
- [15] Keiiti Isii. On sharpness of tchebycheff-type inequalities. *Annals of the Institute of Statistical Mathematics*, 14(1):185–197, 1962.
- [16] Isaac Meilijson and Arthur Nadas. Convex majorization with an application to the length of critical paths. *J. Appl. Probab.*, 16:671–677, 1979.
- [17] Michael Rinehart and Munther A. Dahleh. Shortest path optimization under limited information. In *Proceedings of the 48th IEEE Conference on Decision and Control*, December 2009.
- [18] S.K. Walley, H.H. Tan, and A.M. Viterbi. Shortest path cost distribution in random graphs with positive integer edge costs. In *Twelfth Annual Joint Conference of the IEEE Computer and Communications Societies*, pages 1023–1032, March 1993.
- [19] Gideon Weiss. Stochastic bounds on distributions of optimal value functions with applications to pert, network flows and reliability. *Oper. Res.*, 34(4):595–605, 1986.